

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

This thesis focuses on enhancing the *process comprehension* of an Industrial Control System (ICS) by employing a dynamic black-box analysis approach to analyze the system's physical process. The proposed analysis methodology involves multiple steps and utilizes a custom tool capable of extracting valuable insights about the behavior of the industrial system from both physical process logs and network traffic logs. The tool and methodology will be validated through a case study based on a real-world scenario.

The ultimate goal of this analysis is to gain a deep understanding of the industrial system under examination, with the aim of uncovering as much knowledge as possible. By comprehensively studying the system's behavior and characteristics, potential attackers can develop targeted and covert attack strategies.

Contents

1	Introduction	1
1.1	Contribution	4
1.2	Outline	5
2	Background on Industrial Control Systems	7
2.1	Industrial Control Systems Architecture	9
2.2	Operational Technology Networks	11
2.2.1	Field I/O Devices Layer	11
2.2.2	Controller Network Layer	12
2.2.2.1	Programmable Logic Controllers	12
2.2.2.2	Remote Terminal Units	16
2.2.3	Area Control Layer	17
2.2.3.1	Supervisory Control And Data Acquisition	17
2.2.3.2	Human-Machine Interface	17
2.2.4	Operations/Control Layer	18
2.2.5	Demilitarized Zone	18
2.2.6	Industrial Protocols	19
2.2.6.1	Modbus	20
2.2.6.2	EtherNet/IP	23
2.2.6.3	Common Industrial Protocol (CIP)	25
3	State of the Art	27
3.1	Literature on Process Comprehension	28

3.2	Ceccato et al.'s black-box dynamic analysis for water-tank systems	30
3.2.1	Testbed	32
3.2.2	Scanning of the System and Data Pre-processing	35
3.2.3	Graphs and Statistical Analysis	37
3.2.4	Invariant Inference and Analysis	38
3.2.5	Business Process Mining and Analysis	40
3.2.6	Application	42
3.2.7	Limitations	48
4	Extending and Generalizing Ceccato et al.'s Framework	59
4.1	The Proposed New Framework	60
4.1.1	Framework Structure	63
4.1.2	Python Libraries and External Tools	64
4.2	Analysis Phases	65
4.2.1	A Little Testbed: Stage 1 of iTrust SWaT System	66
4.2.2	Phase 0: Network Analysis	67
4.2.2.1	Extracting Data from PCAP Files	68
4.2.2.2	Network Information	70
4.2.3	Phase 1: Data Pre-processing	73
4.2.3.1	Subsystem Selection	74
4.2.3.2	Dataset Enrichment	77
4.2.3.3	Datasets Merging	84
4.2.3.4	Preliminary Analysis of the Obtained Sub-system	86
4.2.4	Phase 2: Graphs and Statistical Analysis	89
4.2.5	Phase 3: Invariant Inference and Analysis	92
4.2.5.1	Revised Daikon Output	93
4.2.5.2	Types of Analysis	97
4.2.6	Phase 4: Business Process Analysis	102
4.2.6.1	Process Mining of the Physical Process	103
4.2.6.2	Network Data	107
4.2.7	Summary	108

5 Case study: the iTrust SWaT System	111
5.1 Architecture	111
5.1.1 Physical Process	112
5.1.2 Control and Communication Network	114
5.2 Datasets	116
5.2.1 Our Case Study: the 2015 Dataset	117
6 Our Framework at Work on the iTrust SWaT System	121
6.1 Preliminary Operations	122
6.2 Planning the Analysis Strategy	122
6.3 Reverse Engineering of the iTrust SWaT System	124
6.3.1 Reverse Engineering of PLC1 and PLC2	126
6.3.1.1 Pre-processing - Preliminary Analysis . . .	126
6.3.1.2 Graphs and Statistical Analysis	130
6.3.1.3 Invariant Inference and Analysis	133
6.3.1.4 Business Process Mining and Analysis . .	136
6.3.1.5 Properties	139
6.3.2 Reverse Engineering of PLC2 and PLC3	141
6.3.2.1 Pre-processing - Preliminary Analysis . .	141
6.3.2.2 Graphs and Statistical Analysis	145
6.3.2.3 Invariant Inference and Analysis	149
6.3.2.4 Business Process Mining and Analysis . .	152
6.3.2.5 Properties	154
6.3.3 Reverse Engineering of PLC3 and PLC4	155
6.3.3.1 Pre-processing - Preliminary Analysis . .	155
6.3.3.2 Graphs and Statistical Analysis	158
6.3.3.3 Invariant Inference and Analysis	160
6.3.3.4 Business Process Mining and Analysis . .	161
6.3.3.5 Properties	163
7 Conclusion and Future Work	165
List of Figures	169

List of Tables	173
Listings	175
References	177

Introduction

1 **T**HE advent of Industry 4.0 has sparked significant transformations in
2 the realm of *Industrial Control Systems* (ICS). In recent years, there has
3 been a rapid expansion in the interconnectivity and convergence of IT (*In-*
4 *formation Technology*) and OT (*Operational Technology*) systems. While this
5 integration offers undeniable benefits in terms of operational efficiency
6 and effectiveness of ICSs, it has also exposed these systems to the preva-
7 lent vulnerabilities commonly associated with IT environments. Conse-
8 quently, there has been a parallel rise in **cyber-physical attacks**, which
9 originates in the cyber environment and target the physical processes of
10 these systems. These attacks aim to manipulate or disrupt the normal op-
11 eration of ICSs, posing a considerable risk to their integrity and function-
12 ality.

13 *Programmable Logic Controllers* (PLCs), which are essential elements of
14 ICSs, play a vital role in managing electrical equipment like pumps, valves,
15 centrifuges, and more. Serving as a bridge between the cyber and physi-
16 cal realms, PLCs possess a straightforward structure comprising a central
17 processing module (CPU) and additional modules for handling physical
18 inputs and outputs. The user program operates on the CPU and carries
19 out *scan cycles* that involve tasks such as gathering sensor data, execut-
20 ing controller code, and sending commands to actuator devices. Despite
21 modern controllers incorporating security measures to upload authorized

22 firmware, exploiting PLCs can still result in grave consequences for ICSs
23 including physical damage, operational disruptions, environmental haz-
24 ards, and even threats to human safety. These threats can originate from
25 various sources, such as hackers, insider threats, or state-sponsored actors,
26 and can have significant impacts on operational continuity, safety, and fi-
27 nancial losses. The potential impact of an ICS breach necessitates robust
28 security measures.

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

42 Understanding the process operations and functionalities of ICS from
43 an attacker's perspective is a crucial aspect of conducting targeted cyber-
44 attacks. By comprehending how the systems operate, their vulnerabilities,
45 and potential impact, attackers can devise effective strategies to infiltrate
46 and compromise the ICS environment. Here are some key points related
47 to process comprehension of ICS from an attacker's viewpoint:

- 48 • **System Architecture:** attackers aim to understand the architecture
49 of the ICS, including the components involved, such as sensors, ac-
50 tuators, controllers, and human-machine interfaces. This knowledge
51 helps them identify potential entry points and vulnerabilities within
52 the system;

- 53 • **Industrial Protocols:** familiarity with industrial communication pro-
54 tocols, such as Modbus, DNP3, or OPC, is crucial for attackers. Un-
55 derstanding the protocols enables them to analyze the communica-
56 tion between different ICS components and potentially exploit vul-
57 nerabilities or manipulate the data flow;
- 58 • **Operational Processes:** attackers seek to comprehend the operational
59 processes and workflows within the targeted ICS. This includes un-
60 derstanding the sequence of actions, system states, and interactions
61 between various components. Such knowledge enables attackers to
62 identify critical points where disruptions or malicious actions can
63 have significant consequences;
- 64 • **Process Dependencies:** attackers analyze the dependencies between
65 different processes or systems within the ICS. They aim to identify
66 dependencies that, if disrupted or compromised, could have a cas-
67 cading effect on the overall operation. This helps them strategize
68 targeted attacks for maximum impact;
- 69 • **Operational Technology (OT) and Information Technology (IT) Con-
70 vergence:** attackers recognize the convergence of OT and IT systems
71 within modern ICS environments. Understanding how the ICS in-
72 terfaces with external IT networks and devices helps them identify
73 potential attack vectors, such as vulnerable IT systems or insecure
74 connections;
- 75 • **Attack Surface:** by evaluating the attack surface of the ICS, includ-
76 ing network connectivity, remote access options, and third-party in-
77 tegrations, attackers can identify potential entry points and avenues
78 for exploitation.

79 Attackers can acquire a good level of process comprehension of an
80 ICS without employing *Denial of Service* (DoS) techniques through several
81 methods. For instance, they can utilize **probing** techniques, which involve
82 introducing slight perturbations to the system and observing its response

83 and subsequent return to its original state. By analyzing these reactions,
84 attackers can gain valuable insights into the system's behavior.

85 Another technique is the ***Man-in-the-Middle approach***, where attack-
86 ers intercept and analyze network communications to gather information
87 about the system. This allows them to conduct **reconnaissance** of the sys-
88 tem while avoiding detection.

89 ***Reverse engineering*** is another method attackers can employ. By scan-
90 ning memory logs of network-exposed PLCs and analyzing network traf-
91 fic related to the industrial protocols, they can approximate the physical
92 system's model. This technique enables them to derive useful informa-
93 tion without directly intervening in the physical system, as they can work
94 offline using collected logs and analysis tools.

95 1.1 Contribution

96 The original purpose of this thesis was to validate a specific methodol-
97 ogy designed to attain *process comprehension* in an Industrial Control Sys-
98 tem (ICS) when applied to a real-world scenario. This technology utilizes
99 reverse engineering method and adopts a **black box approach**, involving
100 dynamic analysis of both the physical process and network communica-
101 tions within the system. The primary goal of this methodology is to obtain
102 a comprehensive understanding of the industrial process.

103 To accomplish this objective, a framework developed by the authors of
104 the methodology is utilized. This framework incorporates a series of anal-
105 ysis steps that are employed to facilitate the application of the methodol-
106 ogy. The entire framework and methodology are then applied to a spe-
107 cially implemented virtualized testbed, providing a controlled environ-
108 ment for testing and evaluation purposes and to facilitate the process of
109 achieving process comprehension of an Industrial Control System.

110 As the thesis work advanced, it became progressively apparent that
111 both the methodology and tools employed required revision and expan-
112 sion. In response, an extensive revamping of the framework was con-

¹¹³ ducted, aimed at enhancing and expanding its existing features while in-
¹¹⁴ troducing new ones.

¹¹⁵ Therefore, the **primary contribution** of this thesis is to **enhance the**
¹¹⁶ **original methodology** by refining the existing framework, reassessing the
¹¹⁷ approach for each step of the analysis, and incorporating new features.
¹¹⁸ The objective is to achieve a more comprehensive and thorough process
¹¹⁹ comprehension of the industrial system under investigation. This entails
¹²⁰ expanding the methodology to gather a richer set of information, improv-
¹²¹ ing the accuracy of the analysis, and incorporating additional techniques
¹²² to uncover hidden patterns and relationships within the system.

¹²³ Furthermore, the enhanced methodology will be validated through the
¹²⁴ application to a different and larger case study, distinct from the virtu-
¹²⁵ alized testbed used previously. This real-world scenario will provide a
¹²⁶ robust validation of the methodology's effectiveness in a practical set-
¹²⁷ ting, ensuring its applicability and reliability in diverse industrial envi-
¹²⁸ ronments.

¹²⁹ 1.2 Outline

¹³⁰ The thesis is structured as follows:

¹³¹ **Chapter 2:** provides a background on Industrial Control Systems, (ICSS)
¹³² describing their structure, components, and some of the network
¹³³ communication protocols used;

¹³⁴ **Chapter 3:** following an introductory section that provides a brief overview
¹³⁵ of the existing literature on process comprehension in industrial con-
¹³⁶ trol systems, this chapter focuses on a specific paper that outlines
¹³⁷ a methodology to attaining process comprehension of an industrial
¹³⁸ system by employing dynamic blackbox analysis;

¹³⁹ **Chapter 4:** outlines a proposal to improve and extend the methodology
¹⁴⁰ outlined in the previous chapter;

- ¹⁴¹ **Chapter 5:** presents the case study on which the proposed methodology will be applied;
- ¹⁴³ **Chapter 6:** shows how the proposed methodology is applied to the case study illustrated above;
- ¹⁴⁵ **Chapter 7:** outlines final conclusions and future work.

Background on Industrial Control Systems

¹⁴⁶ **I**NDUSTRIAL control systems (ICSs) are information systems used to con-
¹⁴⁷ trol industrial processes such as manufacturing, product handling, pro-
¹⁴⁸ duction, and distribution [1]. ICSs are often found in critical infrastructure
¹⁴⁹ facilities such as power plants, oil and gas refineries, and chemical plant

¹⁵⁰ ICSs are different from traditional IT systems in several key ways. Firstly,
¹⁵¹ ICSs are designed to control physical processes, whereas IT systems are
¹⁵² designed to process and store data. This means that ICSs have different
¹⁵³ requirements for availability, reliability, and performance. Secondly, ICSs
¹⁵⁴ are typically deployed in environments that are harsh and have limited
¹⁵⁵ resources, such as extreme temperatures and limited power. Thirdly, the
¹⁵⁶ protocols hardware and software used in ICSs are often proprietary.

¹⁵⁷ ICSs are becoming increasingly connected to the internet and other net-
¹⁵⁸ works, which has led to increased concerns about their security. Industrial
¹⁵⁹ systems were not originally designed with security in mind, and many of
¹⁶⁰ them have known vulnerabilities that could be exploited by attackers. Ad-
¹⁶¹ ditionally, the use of legacy systems and equipment can make it difficult
¹⁶² to implement security measures. As a result, ICSs are increasingly seen
¹⁶³ as a potential target for cyber attacks, which could have serious conse-
¹⁶⁴ quences for the safe and reliable operation of critical infrastructure: some
¹⁶⁵ notorious examples of cyber attacks are (*i*) the **STUXnet** worm [2], which

¹⁶⁶ purpose was to sabotage the nuclear centrifuges of the enrichment plant
¹⁶⁷ at the Natanz nuclear facility in Iran; (ii) **Industroyer** [3], also referred
¹⁶⁸ as *Crashoverride*, responsible for the attack on the Ukrainian power grid
¹⁶⁹ on December 17, 2016; (iii) the attack on February, 2021 to a water treat-
¹⁷⁰ ment plant in Oldsmar, Florida [4], where the level of sodium hydroxide
¹⁷¹ was intentionally increased to a level approximately 100 times higher than
¹⁷² normal.

¹⁷³ The increasing connectivity of ICSs and the associated security risks
¹⁷⁴ have led to a growing interest in the field of ICS security. Researchers
¹⁷⁵ and practitioners are working to develop new security technologies, stan-
¹⁷⁶ dards, and best practices to protect ICSs from cyber attacks. This includes
¹⁷⁷ efforts to improve the security of ICS networks and devices, as well as the
¹⁷⁸ development of new monitoring and detection techniques to identify and
¹⁷⁹ respond to cyber attacks.

¹⁸⁰ Table 2.1 summarizes the differences between traditional IT and ICSs [5]:

	Traditional IT	ICSs
Focus	Data	Asset
Update Frequency	High	Low
Priority	Confidentiality	Availability
	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)

181 2.1 Industrial Control Systems Architecture

182 In the past, there has been a clear division between *Information Technology* (IT)
183 and *Operational Technology* (OT), both at the technical and organiza-
184 tional levels. Each domain has maintained its own distinct technology
185 stacks, protocols, and standards. However, with the emergence of Indus-
186 try 4.0 and the rapid expansion of industrial automation, which heavily
187 relies on IT tools for monitoring and controlling critical infrastructures,
188 the boundary between IT and OT has started to blur. This trend has paved
189 the way for greater integration between these two domains, thus improv-
190 ing productivity and process quality.

191 General ICS architecture consists in **six levels** each representing a func-
192 tionality: this architecture comprising the OT and IT parts is represented
193 in Figure 2.1 [6][5], according to the *Purdue Enterprise Reference Architecture*
194 (**PERA**), or simply **Purdue Model**:

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

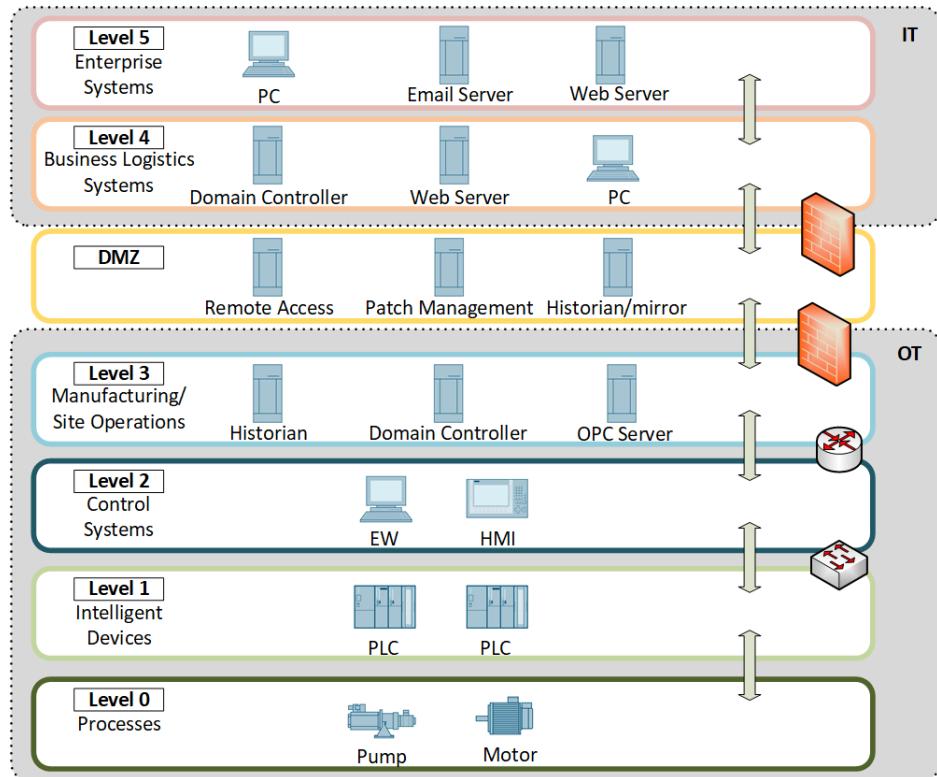


Figure 2.1: ICS architecture schema

- 209 • **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.
- 210
- 211
- 212
- 213
- 214
- 215 • 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.
- 216
- 217
- 218
- 219
- 220 • Level 5 (**Enterprise Systems**): represents the enterprise network, used for the business-to-business activities and for business-to-client pur-
- 221

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

As previously discussed, the gap between IT and OT is steadily narrowing. Nowadays, it is increasingly common to encounter IT elements within the OT realm. For example, desktop PCs are now frequently found in OT environments, and industrial devices are interconnected using standard IT communication protocols like TCP and UDP.

2.2 Operational Technology Networks

Operational Technology primarily encompasses the **tangible aspects** of Industrial Control Systems and directly interfaces with the physical processes of the monitored systems. Its main purpose is to **manage and control the procedures** involved in creating and correcting physical value in various equipment.

This section will focus on the key aspects and components of Operational Technology network, with specific reference to the first four levels of the Purdue model previously seen.

2.2.1 Field I/O Devices Layer

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

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

²⁵⁰ 2.2.2 Controller Network Layer

²⁵¹ *Controller Network* layer includes devices that handle data from and
²⁵² to the *Field I/O Devices* layer. This kind of device is capable of gathering
²⁵³ data from sensors, updating its internal state, and activating actuators (for
²⁵⁴ example opening or close a pump that controls the level of a tank), making
²⁵⁵ decisions based on a customized program, known as its control logic.

²⁵⁶ Commonly found within this layer are *Programmable Logic Controllers*
²⁵⁷ (PLCs) and *Remote Terminal Units* (RTUs): in the upcoming sections, we
²⁵⁸ will examine these elements in detail.

²⁵⁹ 2.2.2.1 Programmable Logic Controllers

²⁶⁰ A *Programmable Logic Controller* (PLC) is a **small and specialized in-**
²⁶¹ **dustrial computer** having the capability of controlling complex industrial
²⁶² and manufacturing processes [7].

²⁶³ Compared to relay systems and personal computers, PLCs are opti-
²⁶⁴ mized for control tasks and industrial environments: they are rugged and
²⁶⁵ designed to withdraw harsh conditions such as dust, vibrations, humid-
²⁶⁶ ity and temperature: they have more reliability than personal computers,
²⁶⁷ which are more prone to crash, and they are more compact a require less
²⁶⁸ maintenance than a relay system.

²⁶⁹ Furthermore, I/O interfaces are already on the controller, so PLCs are
²⁷⁰ easier to expand with additional I/O modules (if in a rack format) to man-
²⁷¹ age more inputs and ouputs, without reconfiguring hardware as in relay
²⁷² systems when a reconfiguration occurs.

²⁷³ PLCs are more *user-friendly*: they are not intended (only) for computer
²⁷⁴ programmers, but designed for engineers with a limited knowledge in
²⁷⁵ programming languages: control program can be entered with a simple
²⁷⁶ and intuitive language based on logic and switching operations instead of
²⁷⁷ a general-purpose programming language (*i.e.* C, C++, ...).

278 **PLC Architecture** The basic hardware architecture of a PLC consists of
279 these elements [8]:

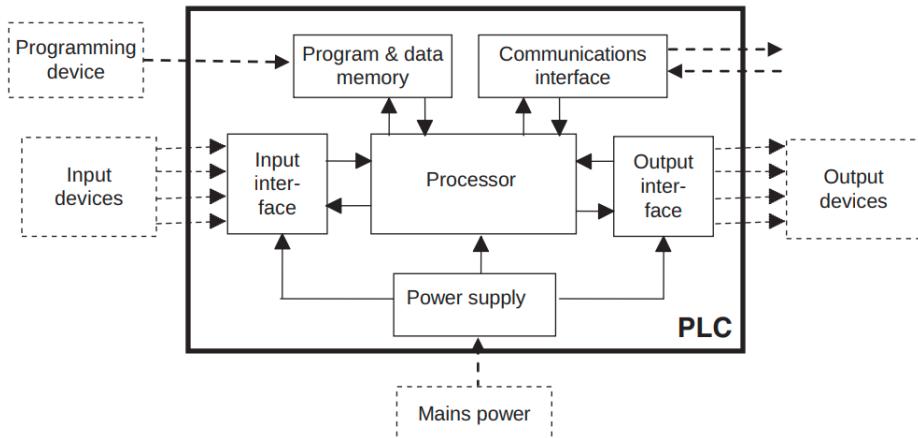


Figure 2.2: PLC architecture

- 280 • **Processor unit (CPU):** contains the microprocessor. This unit inter-
281 pretes the input signals from I/O modules, executes the control pro-
282 gram stored in the Memory Unit and sends the output signals to the
283 I/O Modules. The processor unit also sends data to the Communi-
284 cation interface, for the communication with additional devices.
- 285 • **Power supply unit:** converts AC voltage to low DC voltage.
- 286 • **Programming device:** is used to store the required program into the
287 memory unit.
- 288 • **Memory Unit:** consists in RAM memory and ROM memory. RAM
289 memory is used for storing data from inputs, ROM memory for stor-
290 ing operating system, firmware and user program to be executed by
291 the CPU.
- 292 • **I/O modules:** provide interface between sensors and final control
293 elements (actuators).
- 294 • **Communications interface:** used to send and receive data on a net-
295 work from/to other PLCs.

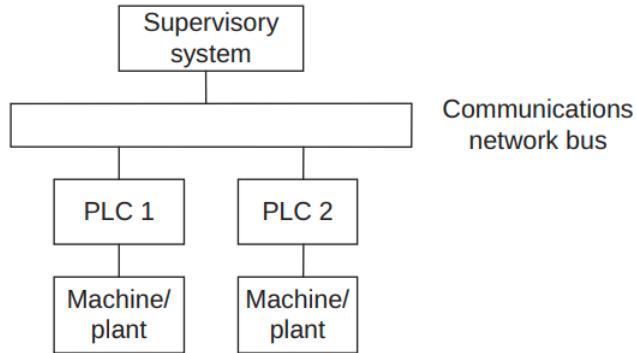


Figure 2.3: PLC communication schema

²⁹⁶ **PLC Programming** Two different programs are executed in a PLC: the
²⁹⁷ **operating system** and the **user program**.

²⁹⁸ The operating system tasks include executing the user program, man-
²⁹⁹ aging memory areas and the *process image table* (memory registers where
³⁰⁰ inputs from sensors and outputs for actuators are stored).

³⁰¹ The user program needs to be uploaded on the PLC via the program-
³⁰² ming device and runs on the process image table in *scan cycles*: each scan
³⁰³ is made up of three phases [9]:

- ³⁰⁴ 1. reading inputs from the process images table
- ³⁰⁵ 2. execution of the control code and computing the physical process evolution
- ³⁰⁷ 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

³¹⁰ Standard PLCs **programming languages** are basically of two types:
³¹¹ **textuals** and **graphicals**. Textual languages include languages such as
³¹² *Instruction List* (IL) and *Structured Text* (ST), while *Ladder Diagrams* (LD),
³¹³ *Function Block Diagram* (FBD) and *Sequential Function Chart* (SFC) belong
³¹⁴ to the graphical languages.

315 Graphical languages are more simple and immediate comparing to the
 316 textual ones and are preferred by programmers because of their features
 317 and simplicity, in particular the **Ladder Logic programming** (see Figure
 318 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 %MMX0..0 : BOOL;
    high AT %MMX0..1 : BOOL;
    open_req AT %MMX0..3 : BOOL;
    close_req AT %MMX0..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

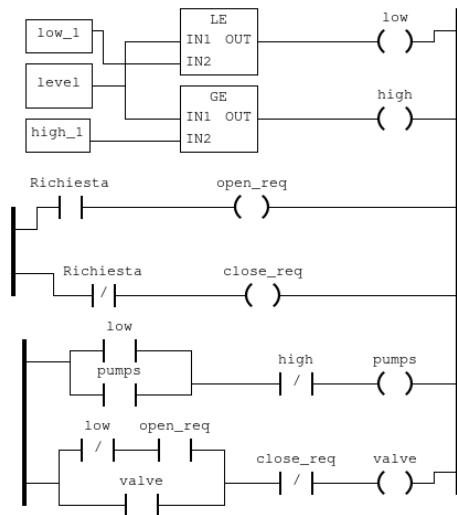
LE3_OUT := LE(level, low_1);
low := LE3_OUT;
GE7_OUT := GE(level, high_1);
high := GE7_OUT;
open_req := Richiesta;
close_req := NOT(Richiesta);
pumps := NOT(high) AND (low OR pumps);
valve := NOT(close_req) AND (open_req AND NOT(low) OR valve);

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

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

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

330 Countermeasures to enhance security in PLC systems may include [10]:

- 331 • protocol modifications implementing **data integrity, authentication**
332 and **protection** against *Replay Attacks*
- 333 • use of *Intrusion Detection and Prevention Systems* (IDP)
- 334 • creation of *Demilitarized Zones* (DMZ) on the network

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

340 2.2.2.2 Remote Terminal Units

341 *Remote Terminal Units* (RTUs) are computers with radio interfacing sim-
342 ilar to PLCs: they transmit telemetry data to the control center or to the
343 PLCs and use messages from the master supervisory system to control
344 connected objects [11].

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

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

³⁵⁹ 2.2.3 Area Control Layer

³⁶⁰ The Area Control layer encompasses hardware and software systems
³⁶¹ useful for supervising, monitoring and controlling the physical process,
³⁶² driving the behavior of the entire infrastructure. The layer includes sys-
³⁶³ tems such as *Supervisory Control and Data Acquisition* (SCADA), *Distributed*
³⁶⁴ *Control Systems* (DCSs), that perform SCADA functions but are usually de-
³⁶⁵ ployed locally, and engineer workstations.

³⁶⁶ 2.2.3.1 Supervisory Control And Data Acquisition

³⁶⁷ *Supervisory Control And Data Acquisition (SCADA)* is a system of soft-
³⁶⁸ ware and hardware elements that allows industrial organizations to [12]:

- ³⁶⁹ • Control industrial processes locally or at remote locations;
- ³⁷⁰ • Monitor, gather, and process real-time data;
- ³⁷¹ • Directly interact with devices such as sensors, valves, pumps, mo-
³⁷² tors, and more through human-machine interface (HMI) software;
- ³⁷³ • Record and aggregate events to send to historian server.

³⁷⁴ The SCADA software processes, distributes, and displays the data, help-
³⁷⁵ ing operators and other employees analyze the data and make important
³⁷⁶ decisions.

³⁷⁷ 2.2.3.2 Human-Machine Interface

³⁷⁸ The *Human-Machine Interface* (HMI) is the hardware and software inter-
³⁷⁹ face that operators use to monitor the processes and interact with the ICS.
³⁸⁰ A HMI shows the operator and authorized users information about sys-
³⁸¹ tem status and history; it also allows them to configure parameters on the
³⁸² ICS such as set points and, send commands and make control decisions
³⁸³ [13].

³⁸⁴ The HMI can be in the form of a physical panel, with buttons and indi-
³⁸⁵ cator lights, or PC software as shown in Figure 2.5.

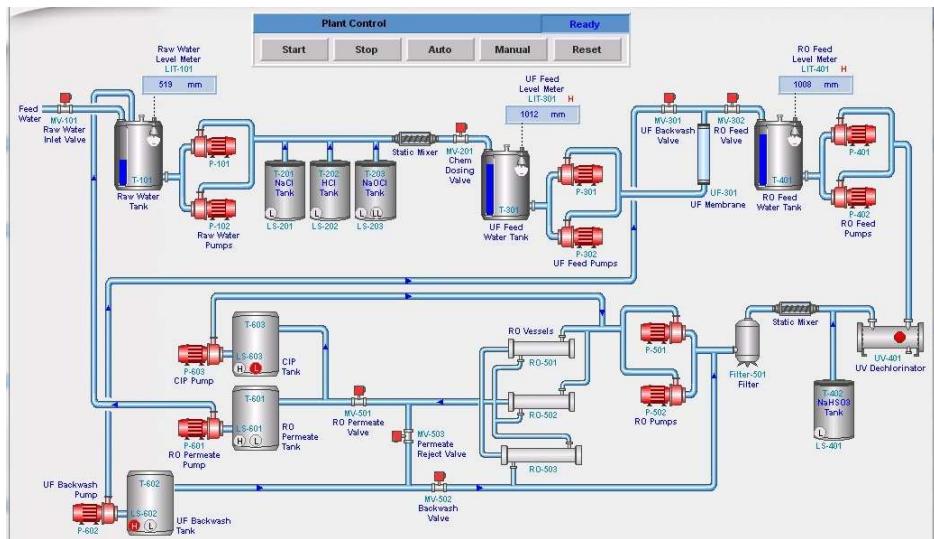


Figure 2.5: Example of HMI for a water treatment plant

386 2.2.4 Operations/Control Layer

387 Within this zone, there are specialized OT devices that are utilized to
 388 manage production workflows on the shop floor [14]. These devices in-
 389 clude:

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

396 2.2.5 Demilitarized Zone

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

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

406 **2.2.6 Industrial Protocols**

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

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

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

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

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

⁴²⁹ **2.2.6.1 Modbus**

⁴³⁰ *Modbus* is a serial communication protocol developed by Modicon (now
⁴³¹ Schneider Electric) in 1979 for use with its PLCs [15] and designed ex-
⁴³² pressly for industrial use: it facilitates interoperability of different devices
⁴³³ connected to the same network (sensors, PLCs, HMIs, ...) and it is also
⁴³⁴ often used to connect RTUs to SCADA acquisition systems.

⁴³⁵ Modbus is the most widely used communication protocol among in-
⁴³⁶ dustrial systems because it has several advantages:

- ⁴³⁷ • simplicity of implementation and debugging
- ⁴³⁸ • it moves raw bits and words, letting the individual vendor to repre-
⁴³⁹ sent the data as it prefers
- ⁴⁴⁰ • it is, nowadays, an **open** and *royalty-free* protocol: there is no need
⁴⁴¹ to sustain licensing costs for implementation and use by industrial
⁴⁴² device vendors

⁴⁴³ Modbus is a **request/response** (or *master/slave*) protocol: this makes it
⁴⁴⁴ independent of the transport layer used.

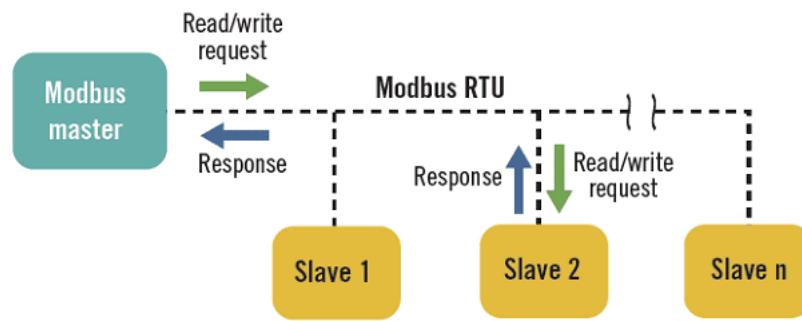


Figure 2.6: Modbus Request/Response schema

⁴⁴⁵ In this kind of architecture, a single device (master) can send requests
⁴⁴⁶ to other devices (slaves), either individually or in broadcast: these slave
⁴⁴⁷ devices (usually peripherals such as actuators) will respond to the master

448 by providing data or performing the action requested by the master using
 449 the Modbus protocol. Slave devices cannot generate requests to the master
 450 [16].

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

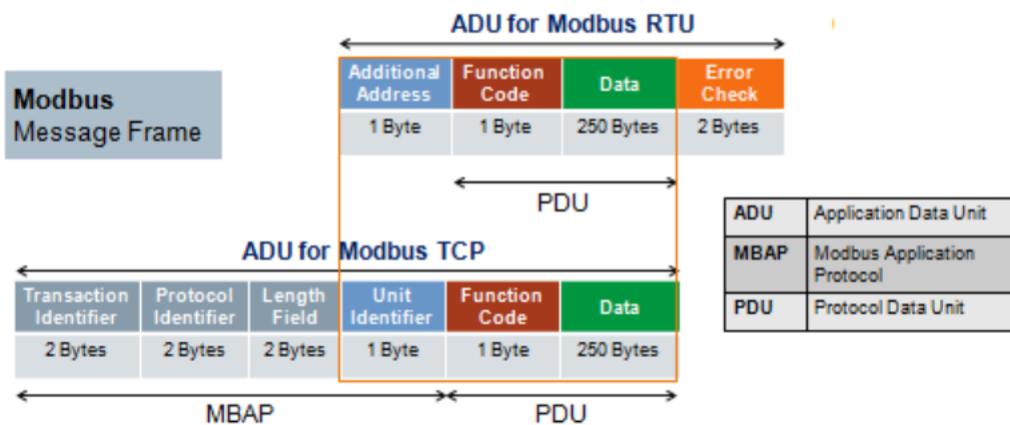


Figure 2.7: Modbus RTU frame and Modbus TCP frame

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

- 458 • *Coil*: binary type, read/write accessible by both masters and slaves
- 459 • *Discrete Input*: binary type, accessible in read-only mode by masters
 460 and in read/write mode by slaves
- 461 • *Analog Input*: 16 bits in size (word), are accessible in read-only mode
 462 by masters and in read/write mode by slaves
- 463 • *Holding Register*: 16 bits in size (word), accessible in read/write mode
 464 by both masters and slaves. Holding Registers are the most com-
 465 monly used registers for output and as general memory registers.

⁴⁶⁶ **Modbus Function Codes** *Modbus Function Codes* are specific codes used
⁴⁶⁷ by the Modbus master within a request frame (see Figure 2.7) to tell the
⁴⁶⁸ Modbus slave device which register type to access and which action to
⁴⁶⁹ perform on it.

⁴⁷⁰ Two types of Function Codes exists: for data access and for diagnostic
⁴⁷¹ 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

⁴⁷² **Modbus Security Issues** Despite its simplicity and widespread use, the
⁴⁷³ Modbus protocol does not have any security feature, which exposes it to
⁴⁷⁴ vulnerabilities and attacks.

⁴⁷⁵ Data in Modbus are transmitted unencrypted (*lack of confidentiality*), with
⁴⁷⁶ no data integrity controls (*lack of integrity*) and authentication checks (*lack*
⁴⁷⁷ *of authentication*), in addition to the *lack of session*. Hence, the protocol is
⁴⁷⁸ vulnerable to a variety of attacks, such as Denial of Services (DoS), buffer
⁴⁷⁹ overflows and reconnaissance activities.

⁴⁸⁰ The easiest attack to bring to the Modbus protocol, however, is **packet**
⁴⁸¹ **sniffing** (Figure 2.8): since, as mentioned earlier, network traffic is un-
⁴⁸² encrypted and the data transmitted is in cleartext, it is sufficient to use
⁴⁸³ a packet sniffer to capture the network traffic, read the packets and thus

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

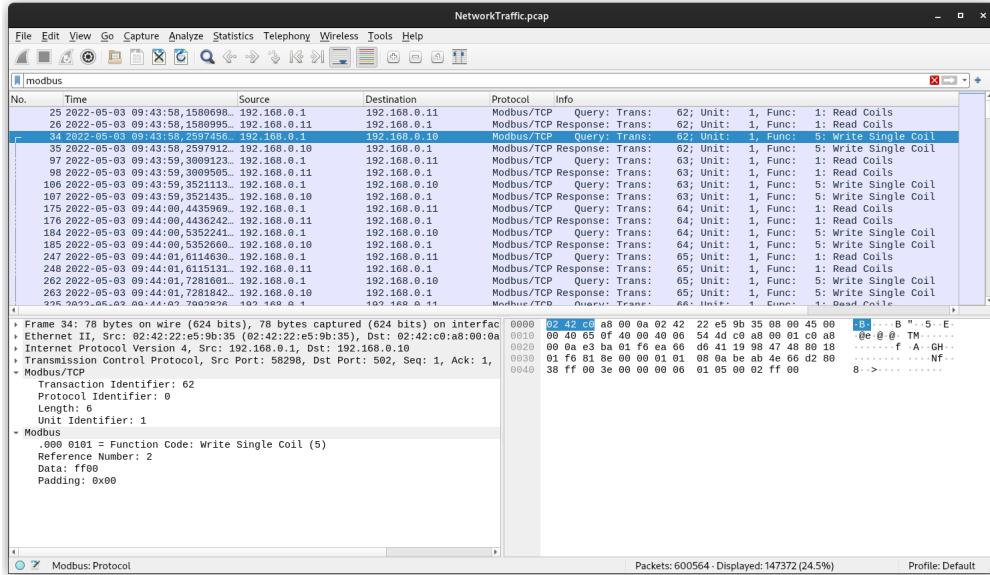


Figure 2.8: Example of packet sniffing on the Modbus protocol

486 To make the Modbus protocol more secure, an encapsulated version
 487 was developed within the *Transport Security Layer* (TLS) cryptographic
 488 protocol, also using mutual authentication. This version of the Modbus
 489 protocol is called **Secure Modbus** or **Modbus TLS**. In addition to this,
 490 Secure Modbus also includes X.509-type certificates to define permissions
 491 and authorisations [17].

492 2.2.6.2 EtherNet/IP

493 *EtherNet/IP* (where IP stands for *Industrial Protocol*) is an open indus-
 494 trial protocol that allows the *Common Industrial Protocol* (CIP) to run on a
 495 typical Ethernet network [18]. It is supported by ODVA [19].

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

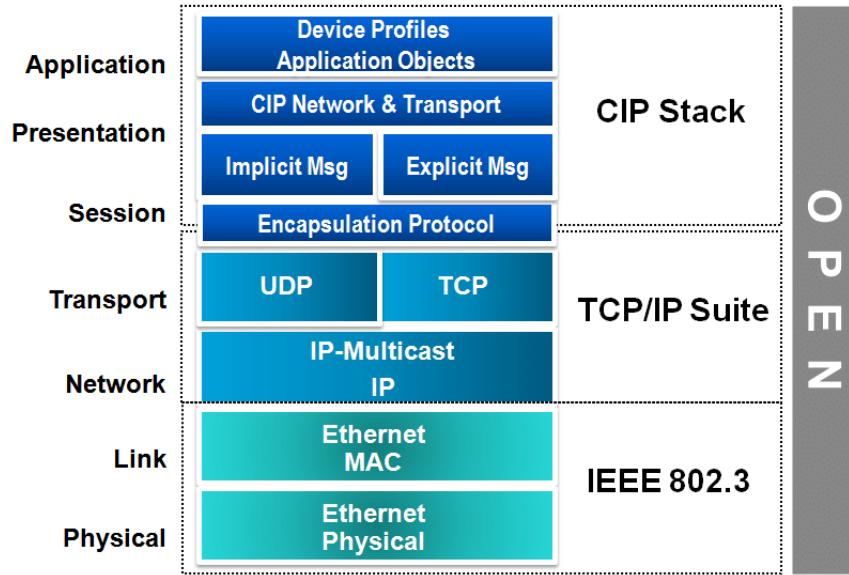


Figure 2.9: OSI model for EtherNet/IP stack

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

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

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

509 **Transport layer** At the transport level, EtherNet/IP encapsulates mes-
510 sages from the CIP stack into an Ethernet message, so that messages can
511 be transmitted from one node to another on the network using the TCP/IP
512 protocol. EtherNet/IP uses two forms of messaging, as defined by CIP
513 standard [18][20]:

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

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

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

522 EtherNet/IP has two types of message connection [18]:

- 523 – **explicit messaging:** *point-to-point* connections to facilitate *request-*
524 *response* transactions between two nodes. These connections use
525 TCP/IP service on port 44818 to transmit messages over Ether-
526 net.
- 527 – **implicit messaging:** this kind of connection moves application-
528 specific **real-time I/O data** at regular intervals. It uses multicast
529 *producer-consumer* model in contrast to the traditional *source-*
530 *destination* model and UDP/IP service (which has lower proto-
531 col overhead and smaller packet size than TCP/IP) on port 2222
532 to transfer data over Ethernet.

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

536 2.2.6.3 Common Industrial Protocol (CIP)

537 The *Common Industrial Protocol* (CIP) is an open industrial automation
538 protocol supported by ODVA. It is a **media independent** (or *transport in-*
539 *dependent*) protocol using a *producer-consumer* communication model and
540 providing a **unified architecture** throughout the manufacturing enterprise
541 [21][22].

542 CIP has been adapted in different types of network:

- 543 • **EtherNet/IP**, adaptation to *Transmission Control Protocol* (TCP) tech-
544 nologies
- 545 • **ControlNet**, adaptation to *Concurrent Time Domain Multiple Access*
546 (CTDMA) technologies
- 547 • **DeviceNet**, adaptation to *Controller Area Network* (CAN) technolo-
548 gies
- 549 • **CompoNet**, adaptation to *Time Division Multiple Access* (TDMA) tech-
550 nologies

551 **CIP objects** CIP is a *strictly object oriented* protocol at the upper layers:
552 each object of CIP has **attributes** (data), **services** (commands), **connec-
553 tions**, and **behaviors** (relationship between values and services of attributes)
554 which are defined in the **CIP object library**. The object library supports
555 many common automation devices and functions, such as analog and dig-
556 ital I/O, valves, motion systems, sensors, and actuators. So if the same
557 object is implemented in two or more devices, it will behave the same way
558 in each device [23].

559 **Security** [24] In EtherNet/IP implementation, security issues are the same
560 as in traditional Ethernet, such as network traffic sniffing and spoofing.
561 The use of the UDP protocol also exposes CIP to transmission route ma-
562 nipulation attacks using the *Internet Group Management Protocol* (IGMP)
563 and malicious traffic injection.

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

State of the Art

567 IN COVENTIONAL IT SYSTEMS, the objective of an attacker is to comprehend
568 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).

574 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.

577 In the context of OT systems, the notion of *reverse engineering* is not limited
578 to its conventional definition, but also includes the concept of **process
579 comprehension**. This term, introduced by Green et al. [25], refers to gaining
580 a comprehensive understanding of the underlying physical process.

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

586 3.1 Literature on Process Comprehension

587 **Keliris and Maniatikos** The first approach presented in this section is by
588 Keliris and Maniatikos [26]: they present a methodology for au-
589 tomating the reverse engineering of ICS binaries based on a *modular*
590 *framework* (called ICSREF) that can reverse binaries compiled with
591 CODESYS [27], one of the most popular and widely used PLC com-
592 pilers, irrespective of the language used.

593 **Yuan et al.** Yuan et al. [28] propose a *data-driven* approach to discover-
594 ing cyber-physical systems process behavior from data directly: to
595 achieve this goal, they have implemented a framework whose pur-
596 pose is to identify physical systems and transition logic inference,
597 and to seek to understand the mechanisms underlying these pro-
598 cesses, making furthermore predictions concerning their state trajec-
599 tories based on the discovered models.

600 **Feng et al.** Feng et al. [29] developed a framework that can generate sys-
601 tem *invariant rules* based on machine learning and data mining tech-
602 niques from ICS operational data log. These invariants are then se-
603 lected by systems engineers to derive IDS systems from them.

604 The experiment results on two different testbeds, the *Water Distri-*
605 *bution system* (WaDi) and the *Secure Water Treatment system* (SWaT),
606 both located at the iTrust - Center for Research in Cyber Security at
607 the University of Singapore of Technology and Design [30], show
608 that under the same false positive rate invariant-based IDSs have a
609 higher efficiency in detecting anomalies than IDS systems based on
610 a residual error-based model.

611 **Pal et al.** Pal et al. [31] work is somewhat related to Feng et al.'s: this
612 paper describes a data-driven approach to identifying invariants au-
613 tomatically using *association rules mining* [32] with the aim of generat-
614 ing invariants sometimes hidden from the design layout. The study
615 has the same objective of Feng et al.'s and uses too the iTrust SwaT

616 System as testbed.

617 Currently this technique is limited to only pair wise sensors and
618 actuators: for more accurate invariants generation, the technique
619 adopted must be capable of deriving valid constraints across multiple
620 sensors and actuators.

621 **Winnicki et al.** Winnicki et al. [33] instead propose a different approach
622 to process comprehension based on the *attacker's perspective* and not
623 limited to mere *Denial of Service* (DoS): their approach is to discover
624 the dynamic behavior of the system, in a semi-automated and process-
625 aware way, through *probing*, that is, slightly perturbing the cyber
626 physical system and observing how it reacts to changes and how
627 it returns to its original state. The difficulty and challenge for the
628 attacker is to perturb the system in such a way as to achieve an ob-
629 servable change, but at the same time avoid this change being seen
630 as a system anomaly by the IDSs.

631 **Green et al.** Green et al. [25] also adopt an approach based on the at-
632 tacker's perspective: this approach consists of two practical exam-
633 ples in a *Man in the Middle* (MitM) scenario to obtain, correlate, and
634 understand all the types of information an attacker might need to
635 plan an attack to alter the process while avoiding detection.

636 The paper shows *step-by-step* how to perform a ICS **reconnaissance**, a
637 phase specifically designed to gather extensive intelligence on mul-
638 tiple fronts, including human factors, network and protocol infor-
639 mation, details about the manufacturing process, industrial applica-
640 tions, and potential vulnerabilities. The primary goal is to accumu-
641 late a wealth of information to enhance understanding and aware-
642 ness in these areas [34]).

643 Reconnaissance phase is fundamental to process comprehension and
644 thus to the execution of MitM attacks.

645 **Ceccato et al.** Ceccato et al. [9] propose a methodology based on a *black*
646 *box dynamic analysis* of an ICS using a reverse engineering tool to

30 3.2 Ceccato et al.'s black-box dynamic analysis for water-tank systems

647 derive from the scans performed on the memory registers of the ex-
648 posed PLCs and the Modbus protocol network scans an approximate
649 model of the physical process. This model is obtained by inferring
650 statistical properties, business process and system invariants from
651 data logs.

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

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

658 3.2 Ceccato et al.'s black-box dynamic analysis 659 for water-tank systems

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

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

672 To accomplish this, the various types of memory registers are analyzed,
673 and attempts are made to determine the nature of the data they might
674 contain.

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

677 To accomplish this, Ceccato et al. have developed a prototype tool [35]
678 that facilitates the reverse engineering of the physical system. This tool
679 goes through four distinct phases:

- 680 1. **scanning of the system and data pre-processing:** this phase involves
681 gathering data to generate data logs for the registers of PLCs and for
682 Modbus network communications.
- 683 2. **graphs and statistical analysis:** The collected data is utilized to pro-
684 vide insights into the memory registers associated with the Modbus
685 protocol by leveraging graphs and statistical data. This analysis ap-
686 proach offers valuable information about the characteristics and pat-
687 terns of the memory registers.
- 688 3. **invariants inference and analysis:** generates system invariants, which
689 are used to identify specific patterns and regularities within the sys-
690 tem. Additionally, this phase provides users with the capability to
691 view invariants related to a particular sensor or actuator.
- 692 4. **business process mining and analysis:** Using event logs, this phase
693 involves reconstructing the business process that depicts how a pro-
694 cess is executed. This step enables a thorough understanding of the
695 sequence of events that occur in the system and how they are in-
696 terrelated, ultimately leading to a comprehensive overview of the
697 business process.

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

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

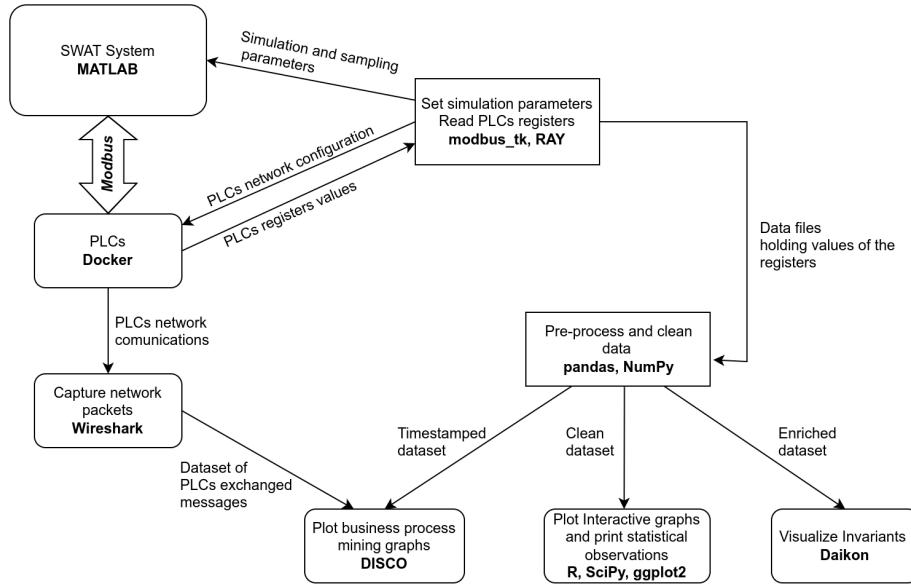


Figure 3.1: Workflow of Ceccato et al.'s stages and operations with used tools

703 3.2.1 Testbed

704 Before delving into the description of the methodology's different phases,
 705 let's first examine the testbed utilized to evaluate this approach. The testbed
 706 employed for testing purposes is a (very) simplified rendition of the iTrust
 707 SWaT system [36], as implemented by Lanotte et al. [37]. Figure 3.2 pro-
 708 vides a graphical representation of the testbed. This simplified version
 709 comprises three stages, each governed by a dedicated PLC.

- 710 **Stage 1** During the initial stage, a **tank** referred to as T-201 with a capacity
 711 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
 712 accumulated water from the tank, directing it to the next stage. Initially,
 713 the water flows from the T-201 tank to the *filtration unit* (which
 714 is not specifically identified by any sensor), and subsequently to a
 715 **second tank** denoted as T-202, with a capacity of 20 gallons.
 716
- 717 **Stage 2** At the second stage, the water stored in tank T-202 flows into the
 718 *reverse osmosis unit* (RO), which serves as both a valve and a continu-

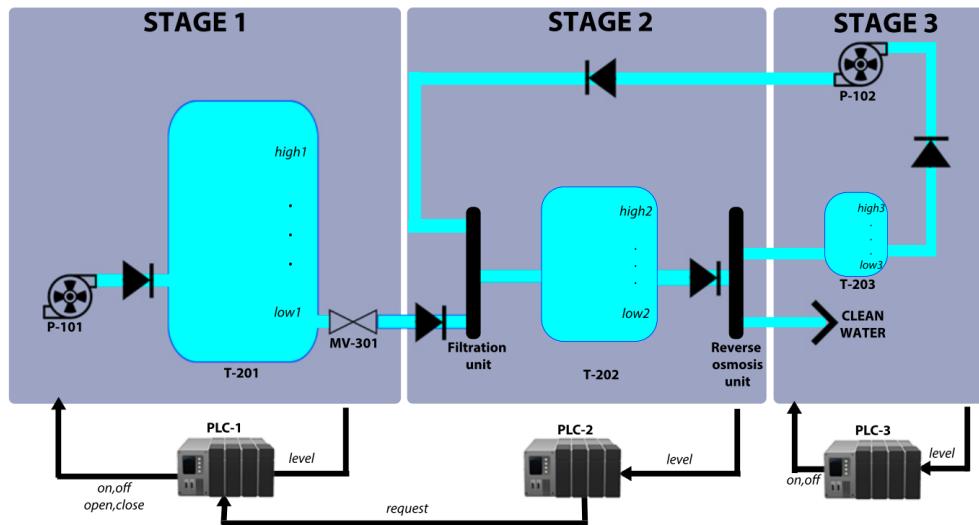


Figure 3.2: The simplified SWaT system used for running Ceccato et al. methodology

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*

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

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

- 739 2. if the level of T-201 reaches the *high setpoint high1* (which is also
740 hardcoded in a specific memory register), then the pump **P-101**
741 **is closed**;
- 742 3. in cases where the level of T-201 is between the *low setpoint low1*
743 and the *high setpoint high1*, PLC1 waits for a request from PLC2
744 to open or close the valve MV-301. If a request to open the valve
745 MV-301 is received, water will flow from T-201 to T-202. How-
746 ever, if no request is received, the valve remains closed. In both
747 situations, the pump P-101 remains closed.

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

- 750 1. when the water level in tank T-202 reaches *low setpoint low2* (also
751 hardcoded in the memory registers), PLC2 sends a request to
752 PLC1 through a Modbus channel to **open valve MV-301.** This
753 request is made in order to allow the water to flow from tank T-
754 201 to tank T-202. The transmission channel between the PLCs
755 is established by copying a boolean value from a memory reg-
756 ister of PLC2 to a corresponding register of PLC1.
- 757 2. when the water level in tank T-202 reaches the *high setpoint high2*
758 value (also hardcoded in the memory registers), PLC2 sends a
759 **close request to PLC1 for valve MV-301.** This request prompts
760 PLC1 to close the valve, stopping the flow of water from tank
761 T-201 to tank T-202.
- 762 3. In cases where the water level in tank T-202 is between the low
763 and high setpoints, the valve MV-301 remains in its current state
764 (open or closed) while the tank is either filling or emptying.

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

- 767 1. If the water level in the backwash tank reaches the *low setpoint*
768 *low3*, **PLC3 sets pump P103 to off**. This allows the backwash
769 tank to be filled.
- 770 2. If the water level in the backwash tank reaches the *high setpoint*
771 *high3*, **PLC3 opens pump P103**. This action triggers the pumping
772 of the entire content of the backwash tank back to the filter
773 unit of T-202.

774 **3.2.2 Scanning of the System and Data Pre-processing**

775 **Scanning tool** The Ceccato et al. scanning tool extends and generalizes
776 a project I did [38] for the "Network Security" and "Cyber Security for IoT"
777 courses taught by Professors Massimo Merro and Mariano Ceccato, re-
778 spectively, in the 2020/21 academic year. The original project involved,
779 in its first part, the recognition within a network of potential PLCs lis-
780 tening on the standard Modbus TCP port 502 using the Nmap module
781 for Python, obtaining the corresponding IP addresses: then a (sequential)
782 scan of a given range of the memory registers of the found PLCs was per-
783 formed to collect the register data. The data thus collected were saved to
784 a file in *JavaScript Object Notation* (JSON) format for later use in the second
785 part of my project.

786 The scanning tool by Ceccato et. al works in a similar way, but extends
787 what originally did by trying to discover other ports on which the Mod-
788 bus protocol might be listening (since in many realities Modbus runs on
789 different ports than the standard one, according to the concept of *security*
790 *by obscurity*) and, most importantly, by **parallelizing and distributing the**
791 **scan** of PLC memory registers through the Ray module [39], specifying
792 moreover the desired granularity of the capture. An example of raw data
793 capture can be seen at Listing 3.1:

```
794 "127.0.0.1/8502/2022-05-03 12_10_00.591": {  
795   "DiscreteInputRegisters": {"%IX0.0": "0"},  
796   "InputRegisters": {"%IW0": "53"},  
797   "HoldingOutputRegisters": {"%QW0": "0"},
```

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

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

Listing 3.1: Example of registers capture

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

803 The tool furthermore offers the possibility, in parallel to the memory
804 registers scan, of **sniffing network traffic** related to the Modbus protocol
805 using the *Man in the Middle* (MitM) technique on the supervisory control
806 network using a Python wrapper for tshark/Wireshark [40] [41]. An ex-
807 ample of raw data obtained with this sniffing can be seen in Listing 3.2:

```
808     Time, Source, Destination, Protocol, Length, Function Code,  
809     ↪ Destination Port, Source Port, Data, Frame length on the  
810     ↪ wire, Bit Value, Request Frame, Reference Number, Info  
811     2022-05-03 11:43:58.158, IP_PLC1, IP_PLC2, Modbus/TCP, 76, Read  
812     ↪ Coils, 46106, 502, , 76, TRUE, 25, , "Response: Trans: 62;  
813     ↪ Unit: 1, Func: 1: Read Coils"
```

Listing 3.2: Example of raw network capture

814 **Data Pre-processing** The data collected by scanning the memory regis-
815 ters of the PLCs are then reprocessed by a Python script and converted
816 in order to create a distinct raw dataset in *Comma Separated Value* for-
817 mat (CSV) for each PLC, containing the memory register values associ-
818 ated with the corresponding controller registers. These datasets are repro-
819 cessed again through the Python modules for **pandas** [42] and **NumPy** [43]
820 by another script to first perform a **data cleanup**, removing all unused reg-
821 isters, **merged** into a single dataset, and finally **enriched** with additional
822 data, such as the **previous value** of all registers and the **measurement**
823 **slope**, that is, the trend of the water level in the system tanks along the
824 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.

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

829 **3.2.3 Graphs and Statistical Analysis**

830 The paper mentions the presence of a *mild graph analysis*, performed
 831 using the framework **R** [44] for statistical analysis at the time of data gath-
 832 ering to find any uncovered patterns, trends and identify measurements
 833 and/or actuator commands through the analysis of registers holding mu-
 834 table values.

835 There is actually no trace of this within the tool: *graph analysis* and *sta-*
 836 *tistical analysis* of the data contained in the PLC memory registers are in-
 837 stead performed using the **matplotlib libraries** and statistical algorithms
 838 made available by the **SciPy libraries** [45], through two separate Python
 839 scripts (see Figure 3.3).

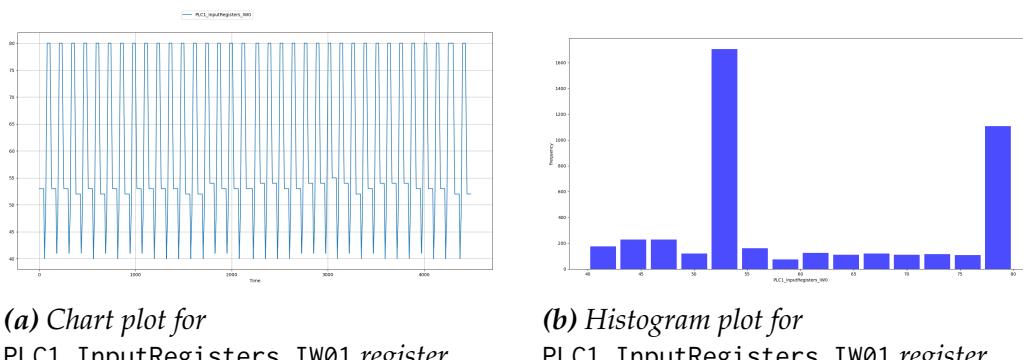


Figure 3.3: Output graphs from graph analysis

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

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

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

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

```
851     Chi-squared test for uniformity
852     Distance      pvalue      Uniform?
853     12488.340    0.00000000    NO
854
855     Shapiro-Wilk test for normality
856     Test statistic   pvalue      Normal?
857     0.844      0.00000000    NO
858
859     Stats of PLC1_InputRegisters_IW0
860     Sample mean = 60.8881; Stddev = 13.0164; max = 80; min =
861     ↪ 40 for 4488 values
```

Listing 3.3: Statistical data for PLC1_InputRegisters_IW0 register

862 3.2.4 Invariant Inference and Analysis

863 For invariant analysis Ceccato et al. rely on **Daikon** [46], a framework
864 to **dynamically detect likely invariants** within a program. An *invariant*
865 is a property that holds at one or more points in a program, properties
866 that are not normally made explicit in the code, but within assert state-
867 ments, documentation and formal specifications: invariants are useful in
868 understanding the behavior of a program (in our case, of the cyber physi-
869 cal system).

870 Daikon uses *machine learning* techniques applied to arbitrary data with
871 the possibility of setting custom conditions for analysis by using a spe-
872 cific file [47] with a *.spinfo* extension (see Listing 3.4). The framework is
873 designed to find the invariants of a program, with various supported pro-
874 gramming languages, starting from the direct execution of the program

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

```
878 PPT_NAME aprogram.point:::POINT
879 VAR1 > VAR2
880 VAR1 == VAR3 && VAR1 != VAR4
```

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

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

```
888 =====
889 aprogram.point:::POINT
890 PLC2_MemoryRegisters_MW1 == PLC3_MemoryRegisters_MW1
891 PLC1_MemoryRegisters_MW0 == 40.0
892 PLC1_MemoryRegisters_MW1 == 80.0
893 PLC1_Coils_QX00 one of { 0.0, 1.0 }
894 [...]
895 =====
896 aprogram.point:::POINT; condition="PLC1_InputRegisters_IW0
897   ↪ > 60"
898 PLC1_InputRegisters_IW0 > PLC1_MemoryRegisters_MW0
899 PLC1_InputRegisters_IW0 > PLC1_Min_safety
900 PLC1_MemoryRegisters_MW0 < prev_PLC1_InputRegisters_IW0
901 [...]
902 =====
903 aprogram.point:::POINT; condition="not(
904   ↪ PLC1_InputRegisters_IW0 > 60)"
905 PLC1_InputRegisters_IW0 < PLC1_MemoryRegisters_MW1
906 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

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

```
907     PLC1_MemoryRegisters_MW1 > prev_PLC1_InputRegisters_IW0  
908     [...]
```

Listing 3.5: The three sections of Daikon analysis outcomes

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

911

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

- 913 • measurements bounded by some setpoint;
- 914 • actuators state changes occurred in the proximity of setpoints or,
915 vice versa, proximity of setpoints upon the occurrence of an actuator
916 state change;
- 917 • state invariants of some actuators correspond to a specific trend in
918 the evolution of the measurements (ascending, descending, or sta-
919 ble) or, vice versa, the measurements trend corresponds to a specific
920 state invariant of some actuators.

921 3.2.5 Business Process Mining and Analysis

922 *Process mining* is the analysis of operational processes based on the
923 event log [48]: the aim of this analysis is to **extract useful informations**
924 from the event data to **reconstruct and understand the behavior** of the
925 business process and how it was actually performed.

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

933 This file is fed to **Disco** [49], a commercial process mining tool, which
934 generates an *activity diagram* similar to UML Activity Diagram and whose

935 nodes represent the activities while the edges represent the relations be-
 936 tween these activities. In Figure 3.4 we can see an example of this diagram
 937 referred to PLC2 of the testbed: nodes represent the trend of register as-
 938 sociated with measurement, actuator state changes, and communications
 939 between PLCs involving these state changes, while edges represent tran-
 940 sitions with their associated time duration and frequency.

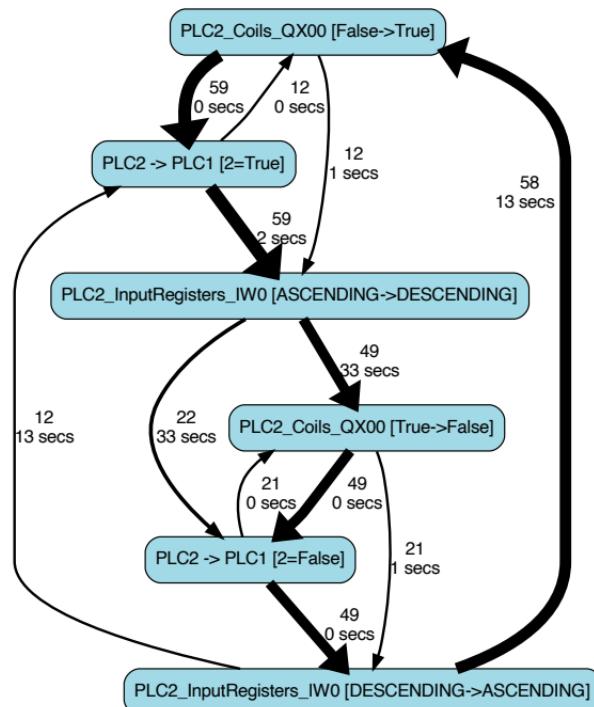


Figure 3.4: An example of Disco generated activity diagram for PLC2

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

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

948 3.2.6 Application

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

957 **Data Collection and Pre-processing** According to what is described in
958 the paper, the data gathering process lasted six hours, with a granular-
959 ity of one data point per second (a full system cycle takes approximately
960 30 minutes). Each datapoint consists of 168 attributes (55 registers plus
961 a special register concerning the tank slope of each PLC) after the en-
962 richment. In addition, IP addresses are automatically replaced by an ab-
963 stract name identified by the prefix PLC followed by a progressive integer
964 (PLC1, PLC2, PLC3), in order to make reading easier.

965 **Graphs and Statistical Analysis** Graphs and Statistical Analysis revealed
966 three properties regarding the contents of the registers:

967 **Property 1:** PLC1_MemoryRegisters_MW0, PLC1_MemoryRegisters_MW1,
968 PLC2_MemoryRegisters_MW0, PLC2_MemoryRegisters_MW1,
969 PLC3_MemoryRegisters_MW0 and PLC3_MemoryRegisters_MW1
970 registers contain constant integer values (40, 80, 10, 20, 0, 10 respec-
971 tively)³. The authors speculate that they may be (relative) hardcoded
972 **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.

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

977 **Property 3:** PLC1_InputRegisters_IW0, PLC2_InputRegisters_IW0 and
 978 PLC3_InputRegisters_IW0 registers contain mutable values.

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

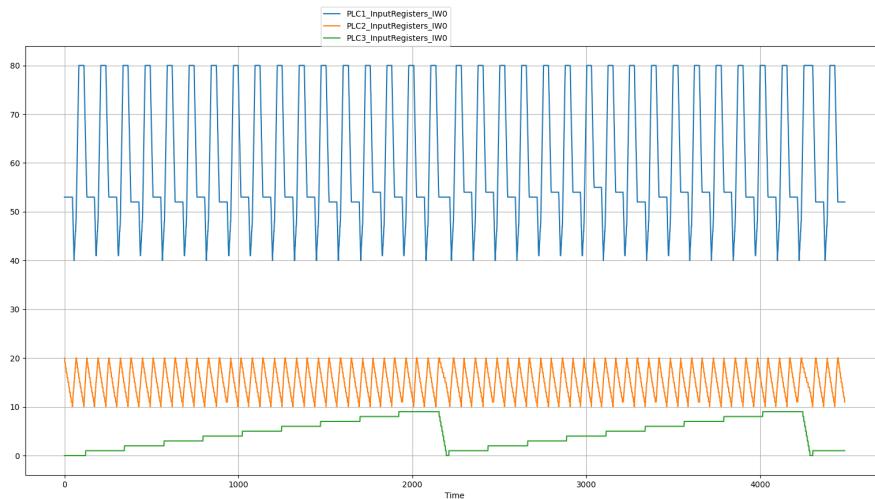


Figure 3.5: Execution traces of *InputRegisters_IW0* on the three PLCs

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

989
 990 **Conjecture 2:**

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

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

- 992 - the relative setpoints for PLC2_InputRegisters_IW0 are 10 and 20;
993 - the relative setpoints for PLC3_InputRegisters_IW0 0 and 9.

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

998 **Business Process Mining and Analysis** With Business Process Mining,
999 the authors aim to **visualize and highlight relevant system behaviors** by
1000 relating PLC states and Modbus commands.

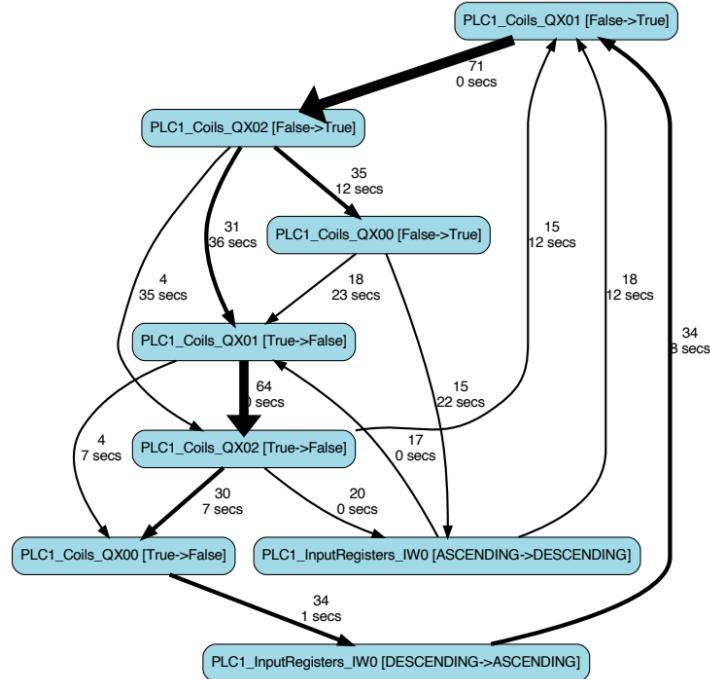
1001 Through analysis of the activity diagrams shown in Figure 3.6, drawn
1002 through Disco, they derive the following properties and conjectures:

1003 **Property 4:** PLC2 sends messages to PLC1 (see Figure 3.6b) which are
1004 recorded to PLC1_Coils_QX02.

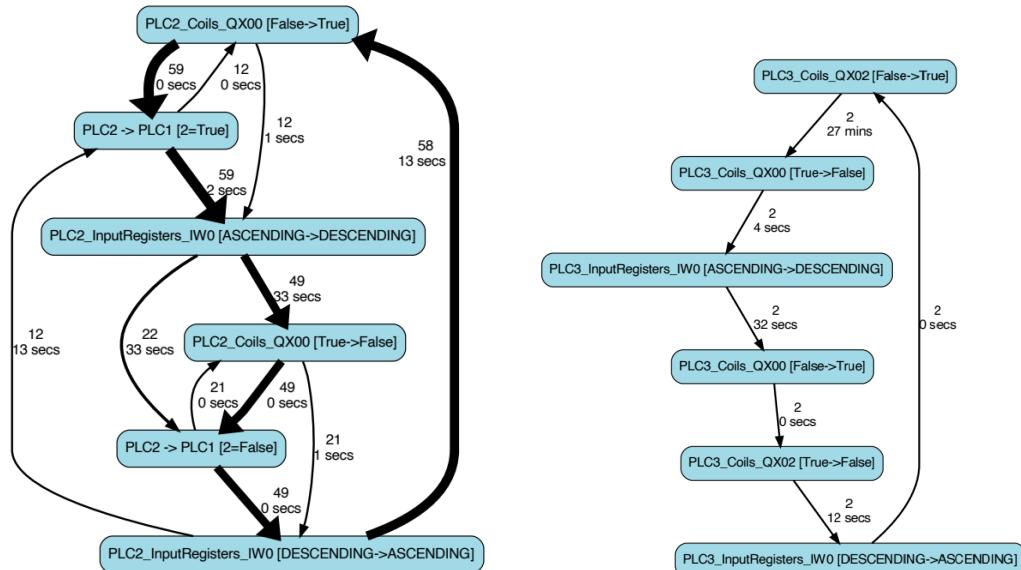
1005 **Conjecture 3:** PLC2_Coils_QX00 determines the trend in tank T-202 (Figure
1006 3.6b).
1007 When this register is set to *True*, the input register PLC2_InputRegisters_IW0
1008 related to the tank controlled by PLC2 starts an **ascending trend**; vice
1009 versa, when the coil register is set to *False*, the input register starts a
1010 **descending trend**.

1011 **Conjecture 4:** If PLC1_Coils_QX00 change his value to True, trend in tank
1012 T-201, related to PLC1_InputRegisters_IW0 and controlled by PLC1,
1013 become **ascending** (see Figure 3.6a)

1014 **Conjecture 5:** PLC3_Coils_QX00 starts a **decreasing trend** in tank T-203, re-
1015 lated to PLC3_InputRegisters_IW0 and controlled by PLC3, whereas
1016 PLC3_Coils_QX02 starts an **increasing trend** on the tank (see Figure
1017 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

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

1018 **Invariant Inference and Analysis** The last phase of the analysis of the
1019 example industrial system is invariant analysis, performed through Daikon
1020 framework. At this stage, an attempt will be made to confirm what has
1021 been seen previously and to derive new properties of the system based on
1022 the results of the Daikon analysis.

1023 To get gradually more and more accurate results, the authors presum-
1024 ably performed more than one analysis with Daikon, including certain
1025 rules within the *splitter information file* (see Section 3.2.4 and Listing 3.4)
1026 based on specific conditions placed on the measurements, for example, the
1027 level of water contained in a tank. Given moreover the massive amount
1028 of invariants generated by Daikon's output, it is not easy to identify and
1029 correlate those that are actually useful for analysis: this must be done man-
1030 ually.

1031 However, it was possible to have confirmation of the conjectures made
1032 in the previous stages of the analysis: starting with the setpoints, analyz-
1033 ing the output of the invariants returned by Daikon⁴ reveals that

1034

```
1035 PLC1_InputRegisters_IW0 >= PLC1_MemoryRegisters_MW0 == 40.0
1036 PLC1_InputRegisters_IW0 <= PLC1_MemoryRegisters_MW1 == 80.0
1037 PLC2_InputRegisters_IW0 >= PLC2_MemoryRegisters_MW0 == 10.0
1038 PLC2_InputRegisters_IW0 <= PLC2_MemoryRegisters_MW1 == 20.0
1039 PLC3_InputRegisters_IW0 >= PLC3_MemoryRegisters_MW0 == 0.0
1040 PLC3_InputRegisters_IW0 <= PLC3_MemoryRegisters_MW1 == 9.0
```

1041 i.e., that the MemoryRegisters_MW0 and MemoryRegisters_MW1 registers of
1042 each PLC contain the **absolute minimum and maximum setpoints**, re-
1043 spectively (*Property 5*).

1045 There is also a confirmation regarding *Property 4*: from the computed
1046 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

1047
1048 PLC1_Coils_QX01 == PLC1_Coils_QX02 == PLC2_Coils_QX00
1049

1050 and from this derive that there is a **communication channel between PLC2**
1051 **and PLC1**, where the value of PLC2_Coils_QX00 is copied to PLC1_Coils_QX01
1052 and PLC1_Coils_QX02 (*Property 6*).

1053 Regarding the **relationships between actuator state changes and mea-**
1054 **surement trends**, invariant analysis yields the results summarized in the
1055 following rules:

1056 **Property 7:** Tank T-202 level *increases* iff PLC1_Coils_QX01 == True. Oth-
1057 erwise, if PLC1_Coils_QX01 == False will be *non-increasing*.

1058 This is because if the coil is *True* the condition

1059 PLC2_InputRegisters_IW0 == PLC2_MemoryRegisters_MW0 == 20.0 && PLC2_slope > 0
1060 is verified. On the opposite hand, if the coil is *False*, the condition
1061 PLC2_InputRegisters_IW0 == PLC2_MemoryRegisters_MW0 == 20.0 && PLC2_slope <= 0 is verified. The
1062 *slope* is increasing if > 0, decreasing if < 0, stable otherwise.

1063 **Property 8:** Tank T-201 level *increases* iff PLC1_Coils_QX00 == True. On the
1064 other hand, if PLC1_Coils_QX00 == False and if PLC1_Coils_QX01 ==
1065 True the level will be *non-decreasing*.

1066 **Property 9:** Tank T-203 level *decreases* iff PLC3_Coils_QX00 == True. It will
1067 be *non-decreasing* if PLC1_Coils_QX00 == False.

1068 The last two properties concern the **relationship between actuator state**
1069 **changes and the setpoints**: it is intended to check what happens to the
1070 actuators when the water level reaches one of these setpoints. From the
1071 analysis of the relevant invariants, the following properties are derived:

1072 **Property 10:** Tank T-201 reaches the upper absolute setpoint when
1073 PLC1_Coils_QX00 changes its state from *True* to *False*. If the coil changes
1074 from *False* to *True*, the tank reaches its absolute lower setpoint.

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

1075 **Property 11:** Tank T-203 reaches the upper absolute setpoint when
1076 PLC3_Coils_QX00 changes its state from *True* to *False*. If the coil changes
1077 from *False* to *True*, the tank reaches its absolute lower setpoint.

1078 3.2.7 Limitations

1079 The methodology proposed by Ceccato et al. is certainly valid and
1080 offers a good starting point for approaching the reverse engineering of
1081 an industrial control system from the attacker's perspective, while also
1082 providing a tool to perform this task.

1083 The limitations of this approach, however, all lie in the tool mentioned
1084 above and also in the testbed described in Section 3.2.1. In this section
1085 we will explain which are the criticisms of each phase, while in Chapter 4
1086 we will formulate proposals to improve and make this methodology more
1087 efficient.

1088 **General Criticism** There are several critical aspects associated with the
1089 application of this approach: the primary one concerns the fact that the
1090 proposed tool seems to be built specifically for the testbed used and that
1091 it is not applicable to other contexts, even to the same type of industrial
1092 control system (water treatment systems, in this case).

1093 What severely limits the analysis performed with the tool implemented
1094 by Ceccato et al. is the use of *ad hoc* solutions and *a posteriori* interventions
1095 done manually on the datasets after the data gathering process: we will
1096 discuss this last aspect in more detail later.

1097 Moreover, there is the presence of many *hardcoded* variables and condi-
1098 tions within the scripts: this makes the system unconfigurable and unable
1099 to properly perform the various stages of the analysis as errors can occur
1100 due to incorrect data and mismatches with the system under analysis.
1101 Having considered, furthermore, only the Modbus protocol for network
1102 communications between the PLCs is another major limiting factor and

1103 does not help the methodology to be adaptable to different systems com-
1104 municating with different protocols (sometimes even multiple ones on the
1105 same system).

1106 Let us now look at the limitations and critical aspects of each phase.

1107 **Testbed** The testbed environment used by Ceccato et. al is entirely simu-
1108 lated, from the physical system to the control system. The PLCs were built
1109 with **OpenPLC** [50] in a Docker environment [51], while the physics part
1110 was built through **Simulink** [52].

1111 OpenPLC is an open source cross-platform software that simulates the
1112 hardware and software functionality of a physical PLC and also offers a
1113 complete editor for PLC program development with support for all stan-
1114 dard languages: *Ladder Logic* (LD), *Function Block Diagram* (FBD), *Instruc-*
1115 *tion List* (IL), *Structured Text* (ST), and *Sequential Function Chart* (SFC).

1116 It is for sure an excellent choice for creating a zero-cost industrial or
1117 home automation and *Internet of Things* (IoT) system that is easy to man-
1118 age via a dedicated, comprehensive and functional web interface. In spite
1119 of these undoubted merits, however, there are (at the moment) **very few**
1120 **supported protocols**: the main one and also referred to in the official doc-
1121 umentation is **Modbus**, while the other protocol is DNP3.

1122 ***First limitation*** The biggest problem with the testbed, however, is not
1123 with the controller part, but with the **physical part**: first of all, it
1124 must be said that although this is something purely demonstrative
1125 even though it is fully functional, the implemented Simulink model
1126 is really **oversimplified** compared to other testbeds. In fact, in the
1127 entire system there are only three actuators, two of which are con-
1128 nected to the same tank and controlled by the same PLC, and sensors
1129 related only to the water level in the system's tanks: in a real sys-
1130 tem there are many more *field devices*, which can monitor and control
1131 other aspects of the system beyond the mere contents of the tanks.
1132 Consider, for example, measuring and controlling the chemicals in

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

1133 the water, the pressure of the liquid in the filter unit, or more simply
1134 the amount of water flow at a given point or time.

1135 All these must be considered and represent a number of additional
1136 variables that makes analysis and consequently reverse engineering
1137 of the system more difficult.

1138 ***Second limitation*** The second critical aspect concerns the **simulation of**
1139 **the physics of the liquid** inside the tanks: Simulink does not con-
1140 sider the fact that inside a tank that is filling (emptying) the liquid
1141 in it undergoes **fluctuations** which cause the level sensor not to see
1142 the water level constantly increasing (decreasing) or at most being
1143 stable at each point of detection. Figure 3.7 exemplifies more clearly
1144 with an example the concept just expressed: these oscillations cause
1145 a **perturbation** in the data.

1146 This issue leads to the difficulty, on a real physical system, of **cor-**
1147 **rectly calculating the trend of a measurement** by using the slope
1148 attribute: if this was obtained with a too low granularity, the trend
1149 will be oscillating between increasing and decreasing even when in
1150 reality this would be in general increasing (decreasing) or stable; on
1151 the other hand, if the slope was obtained with a too high granularity
1152 there is a loss of information and the trend may be "flattened" with
1153 respect to reality.

1154 In the present case, the slope in the Simulink model was calculated
1155 statically with a (very) low granularity, 5 and 6 seconds according
1156 to the Properties 7 and 9 described in the original paper: an aver-
1157 agely careful reader will have already guessed that this granularity
1158 is inapplicable to the real system in Figure 3.7b. As we will later see,
1159 we need to **operate on the data perturbations** to be able to obtain a
1160 suitable granularity and a correct calculation of the slope and conse-
1161 quently of the measurement trend.

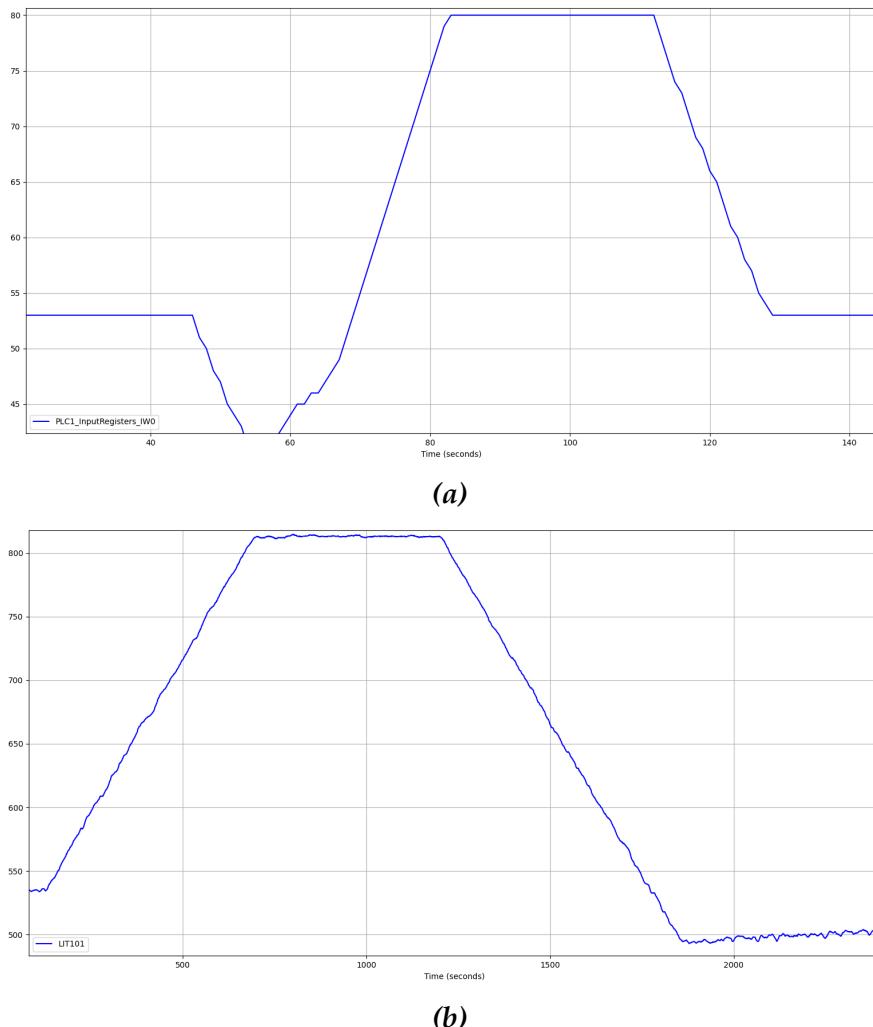


Figure 3.7: Water physics compared: simulated physics in the Simulink model (a) and physics in a real system (iTrust SWaT) (b). Fluctuations in the tank level in (b), almost completely absent in (a), can be appreciated.

1162 **Pre-processing** In the pre-processing phase, the authors make use of a
 1163 Python script to merge all the datasets of the individual PLCs into a single
 1164 dataset, remove the (supposedly) unused registers, and finally enrich the
 1165 obtained dataset with additional attributes. These attributes, as seen in
 1166 3.2.2, are:

- 1167 • the **previous value** of all registers;

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

- 1168 • some **additional relative setpoints** named PLC x _Max_safety and
1169 PLC x _Min_safety (where x is the PLC number), which represent a
1170 kind of alert on reaching the maximum and minimum water levels
1171 of the tanks;
- 1172 • the **measurement slope**.

1173 ***First limitation*** Merging the datasets of all individual PLCs into a single
1174 dataset representing the entire system can be a sound practice if the
1175 system to be analyzed is (very) small as is the testbed analyzed here,
1176 consisting of a few PLCs and especially a few registers. If, however,
1177 the complexity of the system increases, this type of merging can be-
1178 come counterproductive and make it difficult to analyze and under-
1179 stand the data obtained in subsequent steps.

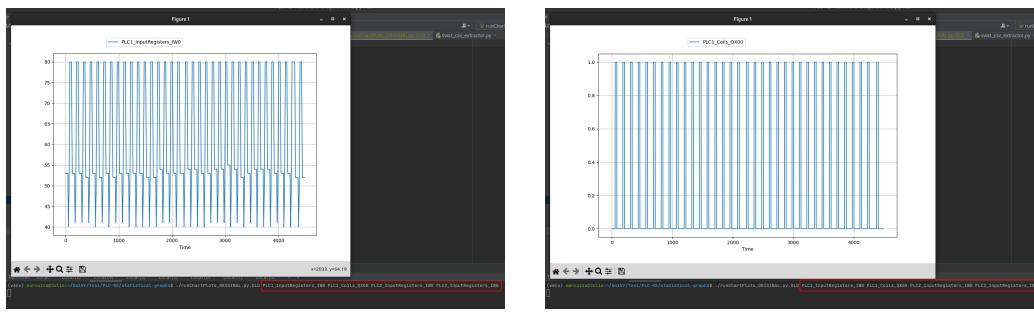
1180 In short, there is no possibility to analyze only a subsystem and thus
1181 make the analysis faster and more understandable. Moreover, a data
1182 gathering can take up to days, and the analyst/attacker may need to
1183 make an analysis of the system isolating precise time ranges, ignor-
1184 ing everything that happens before and/or after: all of this, with the
1185 tool we have seen, cannot be done.

1186 ***Second limitation*** Regarding the additional attributes, looking at the code
1187 of the script that performs the enrichment, we observed that **some at-**
1188 **tributes were manually inserted** after the merging phase: we are re-
1189 ferring in particular to the attributes PLC x _Max_safety and PLC x _Min_safety,
1190 whose references were moreover hardcoded into the script, and the
1191 *slope* whose calculation method we mentioned in the previous para-
1192 graph about the testbed limitations.

1193 In the end, only the attribute *prev* related to the value at the previous
1194 point of the detection is inserted automatically for all registers, more-
1195 over without the possibility to choose whether this attribute should
1196 be extended to all registers or only to a part.

1197 **Graphs and Statistical Analysis** Describing the behavior of graphical
1198 analysis in Section 3.2.3 we had already mentioned that only one register

1199 plot at a time was shown and not, for example, a single window containing
 1200 the charts of all registers entered by the user as input from the com-
 1201 mand line, such as in Figure 3.5. Figure 3.8 shows the actual behavior
 1202 of graphical analysis: note that although we have specified four registers
 1203 (highlighted in red in the figures) as command-line parameters, only one
 1204 at a time is shown and it is necessary to close the current chart in order to
 1205 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

1206 **First limitation** While displaying charts for individual registers still pro-
 1207 vides useful information about the system such as the distinction
 1208 between actuators and measurements and the general trend of the
 1209 latter, single display does not allow one to catch, or at least makes it
 1210 difficult, the relationship that exists between actuators and measure-
 1211 ments, where it exists, because a view of the system as a whole is
 1212 missing.

1213 In this way, the risk is to make conjectures about the behavior of the
 1214 system that may prove to be at least imprecise, if not inaccurate.

1215 **Second limitation** On the other hand, regarding the statistical analysis,
 1216 two observations need to be made: the first is that for the given sys-
 1217 tem, I personally was unable to appreciate the usefulness of the gen-
 1218 erated histogram in Figure 3.3b, as it does not provide any particular
 1219 new information that has not already been obtained from the graph-

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

1220 ical analysis (except maybe something marginal); the second obser-
1221 vation pertains to the presentation of statistical information obtained
1222 from the histogram plot. In certain cases, the histogram plot itself
1223 can overshadow the displayed statistical information. These statis-
1224 tics are actually shown on the terminal from which the script is exe-
1225 cuted. However, to an inattentive or unfamiliar user, these statistics
1226 may be mistaken for debugging output or warnings, as they coin-
1227 cide with the display of the histogram plot window, which takes the
1228 focus (see Figure 3.9).

1229

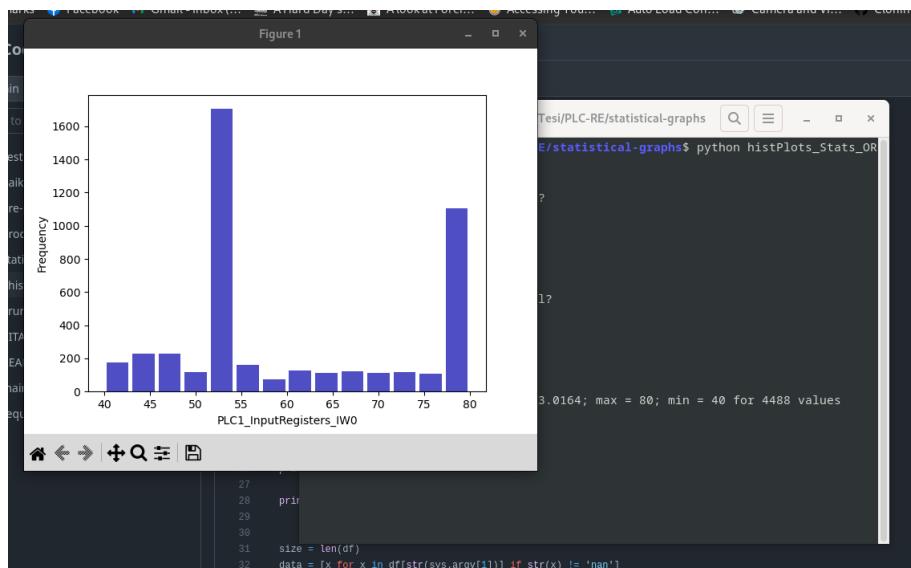


Figure 3.9: Histogram plot overshadowing statistical information shown on the terminal window in the background

1230

In general, however, little statistical information is provided.

1231 **Business Process Mining and Analysis** Concerning the data mining,
1232 this is a purely *ad hoc solution*, designed to work under special conditions:
1233 first, the timestamped dataset of the physical process and the one obtained
1234 after the packet sniffing operation of Modbus traffic on the network need
1235 to be synchronized and have the same granularity, in this case one event
1236 per second.

1237 It is relatively easy, therefore, to find correspondences between Mod-
1238 bus commands sent over the network and events occurring on the physical
1239 system, such as state changes in actuations, due in part to the fact that the
1240 number of communications over the network is really small (see Section
1241 3.2.1).

1242 ***First limitation*** In a real system, network communications are much more
1243 numerous and involve many more devices even in the same second:
1244 finding the exact correspondence with what is happening in the cy-
1245 ber physical system becomes much more difficult.

1246 Since this is, as mentioned, an *ad hoc* solution, only the Modbus pro-
1247 tocol is being considered: as widely used as this industrial protocol
1248 is, other protocols that are widely used [53] such as EtherNet/IP (see
1249 Section 2.2.6.2) or Profinet should be considered in order to extend
1250 the analysis to other industrial systems that use a different commu-
1251 nication network.

1252 ***Second limitation*** The other limiting aspect of the business process min-
1253 ing phase is the **process mining software** used to generate the ac-
1254 tivity diagram. As mentioned in Section 3.2.5, the process mining
1255 software used by Ceccato et al. is **Disco**: this is commercial soft-
1256 ware, with an academic license lasting only 30 days (although free of
1257 charge), released for Windows and MacOS operating systems only,
1258 which makes its use under Linux systems impossible except by us-
1259 ing emulation environments such as Wine.

1260 For what is my personal vision and training as a computer scientist,
1261 it would have been preferable to use a *cross-platform, freely licensed*
1262 *open source* software alternative to Disco: one such software could
1263 have been **ProM Tools** [54], a framework for process mining very
1264 similar to Disco in functionality, but fitting the criteria just described,
1265 or use Python libraries such as **PM4PY** [55], which offer ready-to-use
1266 algorithms suitable for various process mining needs.

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

1267 **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

1273 **First limitation** The biggest problem with Daikon applied to the computation of invariants of an industrial system is the difficult reading of the resulting output: the software in fact returns a very long list 1274 of invariants, one invariant per line, many of no use and without 1275 correlating invariants that may have common features or deriving 1276 1277

1278 additional information from them. The process of screening and rec-
1279ognizing the significant invariants, as well as the correlation between
1280them, must be done by a human: certainly not an easy task given the
1281volume of invariants one could theoretically be faced with (hundreds
1282and hundreds of invariants). An example of Daikon's output can be
1283seen in Figure 3.10.

1284 **Second limitation** The bash script used in this phase of the analysis does
1285not help at all in deriving significant invariants more easily: it merely
1286launches Daikon and saves its output to a text file by simply redirect-
1287ing the stdout to file. No data reprocessing is done during this step.
1288 In addition, if a condition is to be specified to Daikon before perform-
1289ing the analysis, it is necessary each time to edit the .spinfo file by
1290manually entering the desired rule, an inconvenient operation when
1291multiple analyses are to be performed with different conditions each
1292time.

1293 Table 3.1 provides a summary of the limitations discussed regarding the
1294 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

58 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

1295 IN Chapter 3, we presented the state of the art of *process comprehension*
1296 of an Industrial Control System (ICS) with a focus on the methodology
1297 proposed by Ceccato et al. [9][Section 3.2], explaining what it consists of,
1298 its practical application on a testbed, and most importantly highlighting
1299 its limitations and critical issues (see Section 3.2.7).

1300 In this chapter we will present a **proposal to improve the methodology**
1301 seen in the previous chapter, addressing most of the critical issues (or
1302 at least trying to do so) described in Section 3.2.7 by almost completely
1303 rewriting the original framework, enhancing its functionalities and in-
1304 inserting new ones where possible, while preserving its general structure
1305 and approach. Indeed, the system analysis will encompass the same **four**
1306 **phases** as the original methodology (Data Pre-processing, Graphs and Sta-
1307 tistical Analysis, Invariant Inference and Analysis, and Business Process
1308 Mining and Analysis). In addition to these phases, a **fifth phase** dedicated
1309 exclusively to **network traffic analysis** will be introduced. Each phase of
1310 the original methodology will undergo thorough revision to enhance the
1311 understanding of the industrial system being analyzed and its behavior,
1312 aiming to provide a more complete and coherent process comprehension.
1313 This revision aims to enrich the analysis process, making it more robust

¹³¹⁴ and transparent.

¹³¹⁵ Please note that our proposals do not include improvements to the data
¹³¹⁶ gathering phase. This decision is solely because the new framework will
¹³¹⁷ not be tested on the same case study used by Ceccato et al. (Section 3.2.1),
¹³¹⁸ but rather on a different case study known as the iTrust SWaT system [36].
¹³¹⁹ iTrust has already provided some datasets that contain the execution trace
¹³²⁰ of the physical system and the network traffic scan for this case study. In
¹³²¹ contrast to the case study conducted by Ceccato et al., the iTrust SWaT
¹³²² system is not simulated but rather a **real-world system**. This distinction
¹³²³ is important as it introduces additional complexities and considerations
¹³²⁴ when analyzing and interpreting the data.

¹³²⁵ **Outline** The upcoming sections provide an outline of what to expect:

- ¹³²⁶ • **Section 4.1** will introduce a **novel framework** that we have devel-
¹³²⁷ oped to address the limitations of Ceccato et al.'s framework. Sec-
¹³²⁸ tion will provide a brief overview of the framework's features and
¹³²⁹ structure;
- ¹³³⁰ • in **Section 4.2** the **features of our framework** will be demonstrated
¹³³¹ through their application to the different steps of the methodology.
¹³³² Practical examples will be used to illustrate the functionality of the
¹³³³ framework. To facilitate this, we will utilize a more complex and
¹³³⁴ detailed testbed compared to the one employed by Ceccato et al.

¹³³⁵ 4.1 The Proposed New Framework

¹³³⁶ In our version of the framework we decided to follow a few design
¹³³⁷ choices:

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

1340 3. It must provide greater **flexibility and ease of use** at every stage.

1341 In the following, we discuss these three features in more detail.

1342 **Single Programming Language** The implementation of Ceccato et al.'s tool
1343 involved the use of different programming languages for each of the
1344 phases, ranging from Python to Java, and even including Bash script-
1345 ing.

1346 However, we believe that this heterogeneity introduces challenges
1347 in terms of user-friendliness and intuitiveness. It becomes more dif-
1348 ficult for users to navigate and operate the tool effectively. Fur-
1349 thermore, the utilization of multiple technologies complicates code
1350 maintenance and the addition of new features, especially when man-
1351 aged by a single person who may have expertise in only one lan-
1352 guage while being less familiar with others.

1353 Considering these factors, we have made the decision to adopt a
1354 single programming language for the framework to ensure homo-
1355 geneity, simplify usability, and facilitate code maintenance for fu-
1356 ture users. Python has been chosen as the language of choice due
1357 to its simplicity, readability, versatility, and powerfulness. Addition-
1358 ally, Python benefits from a vast ecosystem of libraries and packages
1359 that cater to various requirements, making it a suitable choice for our
1360 framework.

1361 **System Independence** One of the significant limitations we identified in
1362 Ceccato et al.'s tool, as highlighted in Section 3.2.7, is its **strong de-**
1363 **pendence on the specific testbed** being used. This means that the
1364 tool lacks the flexibility to configure parameters for analyzing differ-
1365 ent industrial systems.

1366 To address this limitation and achieve system independence, we have
1367 made a crucial improvement in our framework. We have introduced
1368 a comprehensive configuration file called *config.ini* that allows users
1369 to customize all the necessary parameters for analyzing the targeted

1370 system. This general configuration file eliminates any references to
1371 hardcoded variables or values found in the Ceccato et al.'s tool, pro-
1372 viding users with the flexibility to tailor the analysis according to
1373 their specific requirements.

1374 **Flexibility and Ease on Use** The lack of flexibility and ease of use in a tool
1375 can be a significant disadvantage, as it hampers its effectiveness and
1376 makes it difficult for users to achieve their desired outcomes. Cec-
1377 cato et al.'s tool was affected by these limitations, as it required users
1378 to run scripts from the command line without sufficient options or
1379 parameters for customizing the analysis. Consequently, the tool fell
1380 short in terms of user-friendliness and failed to provide the neces-
1381 sary flexibility to accommodate specific user requirements.

1382 To settle these issues, we enhanced the command-line interface in the
1383 proposed framework by adding new options and parameters. These
1384 new features provide the user with greater flexibility, enabling to
1385 specify parameters and options that allow for more in-depth anal-
1386 ysis and focused results analyzing data more effectively and effi-
1387 ciently. With these enhancements, the framework has become more
1388 user-friendly, reducing the learning curve and making it more acces-
1389 sible to a wider range of users.

1390 This, in turn, makes the framework more valuable and useful, in-
1391 creasing its adoption and effectiveness across a range of industrial
1392 control systems and applications.

1393 Moreover, with new options and parameters users can now access a
1394 range of customizable options and parameters, making the tool more
1395 intuitive and user-friendly.

1396 Overall, the enhancements made to the framework represent a signifi-
1397 cant step forward in making it more effective, efficient, and user-friendly.

1398 4.1.1 Framework Structure

1399 The proposed framework follows a similar structure to the original
 1400 tool, with a division into five main directories representing different phases
 1401 of the analysis: **data pre-processing, graphs and statistical analysis, in-**
1402 variant analysis, and process mining. A new phase is added compared to
 1403 the original, concerning the **analysis of the network traffic**. These directo-
 1404 ries contain the corresponding Python scripts responsible for performing
 1405 the analysis, along with any necessary subdirectories and input/output
 1406 files to ensure the proper functioning of the framework.

```
1407 .
1408   |-- config.ini
1409   |-- daikon
1410     |-- Daikon_Invariants
1411     |-- daikonAnalysis.py
1412     |-- runDaikon.py
1413   |-- network-analysis
1414     |-- data
1415     |-- networkAnalysis.py
1416     |-- export_pcap_data.py
1417     |-- swat_csv_extractor.py
1418   |-- pre-processing
1419     |-- mergeDatasets.py
1420     |-- system_info.py
1421   |-- process-mining
1422     |-- data
1423     |-- process_mining.py
1424   |-- statistical-graphs
1425     |-- histPlots_Stats.py
1426     |-- runChartSubPlots.py
```

1427 *Listing 4.1: Novel Framework structure and Python scripts*

1427 The *config.ini* file is located in the root directory of the framework. This
 1428 file holds significant importance as it grants the framework independence
 1429 from the industrial control system being analyzed. Within this file, users
 1430 have the opportunity to configure various general parameters and op-
 1431 tions. These include specifying file paths for reading or writing files, as

1432 well as options related to specific analysis phases.

1433 The file is organized into sections, with each section dedicated to a spe-
1434 cific aspect of the configuration. These sections contain user-customizable
1435 parameters that are later referenced within the Python scripts comprising
1436 the framework. Sections of *config.ini* are:

- 1437 • [PATHS]: defines general paths such as the project root directory;
- 1438 • [PREPROC]: includes parameters needed for the **pre-processing phase**,
1439 like the directory containing the raw datasets of the individual PLCs
1440 and *granularity* for the slope calculation. The granulariy is given in
1441 terms of the time interval over which the slope is calculated;
- 1442 • [DATASET]: defines settings and parameters used during the **dataset**
1443 **enrichment** stage, for example the additional attributes;
- 1444 • [DAIKON]: defines parameters needed for **invariant analysis** with
1445 Daikon, e.g. directories and files containing the outcomes of the anal-
1446 ysis;
- 1447 • [MINING]: contains parameters used during the **process mining**
1448 phase, such as data directory;
- 1449 • [NETWORK]: includes specific settings for extracting the data ob-
1450 tained from the packet sniffing phase on the ICS network and con-
1451 verting it to CSV format. It also defines the **network protocols** that
1452 are to be analyzed.

1453 4.1.2 Python Libraries and External Tools

1454 As the framework has been developed entirely in Python, the objec-
1455 tive was to minimize reliance on external tools and instead integrate vari-
1456 ous functionalities within the framework itself. The aim was to make the
1457 framework independent from external software. The only remaining ex-
1458 ternal tool from the Ceccato et al. tool is Daikon. This choice was made

1459 because there is currently no better alternative or Python package avail-
1460 able that offers the same functionalities as Daikon.

1461 Instead, the framework extensively utilizes Python libraries for han-
1462 dling various functionalities and input data. The core libraries on which
1463 the framework relies are:

- 1464 • **Pandas**, also used in the Ceccato et al.'s tool for dataset management,
1465 but whose use here has been deepened and extended
- 1466 • **NumPy**, often used together with Pandas to perform some opera-
1467 tions to support it;
- 1468 • **MatPlotLib**, for managing and plotting graphical analysis;
- 1469 • other scientific libraries such as **SciPy**, **StatsModel** [56] and **Net-**
1470 **workX** [57], for mathematical, statistical and analysis operations on
1471 the data;
- 1472 • **GraphViz**, for the creation of activity diagrams in the process mining
1473 phase.

1474 Having now seen the structure of the framework, in the next sections we
1475 will go into more detail describing our proposal.

1476 4.2 Analysis Phases

1477 In this section, the behavior of the proposed framework throughout the
1478 various analysis phases of the methodology will be illustrated. Here is a
1479 brief **outline** of the content covered in this section:

- 1480 • **Section 4.2.1:** this section will present the **testbed** used to demon-
1481 strate examples of the framework's application;
- 1482 • **Section 4.2.2:** the introduction of the new **network traffic analysis**
1483 phase (*Phase 0*) will be discussed. This phase aims to understand the
1484 structure of the industrial system's network and analyze its commu-
1485 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.

1501 4.2.1 A Little Testbed: Stage 1 of iTrust SWaT System

Before we proceed with presenting the analysis steps of the proposed framework, let us introduce the testbed that will serve as an illustration for practical examples, demonstrating the effectiveness of our methodology and the potential of the framework. This testbed corresponds to Stage 1 of the iTrust SWaT (Secure Water Treatment) system [36]. The selection of this testbed is intentional, as it serves as a precursor to the comprehensive case study we will be addressing in the upcoming chapters. The iTrust SWaT system offers elements of greater complexity compared to the individual stages of the Ceccato et al. testbed.

1511 The testbed comprises several components, including:

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

- 1513 • a PLC responsible for monitoring and controlling the operations within
1514 the stage;
- 1515 • **two sensors** that provide readings of the water level within the tank
1516 and the incoming flow. These sensors are identified as LIT101 and
1517 FIT101 in the PLC registers;
- 1518 • **three actuators**, namely a valve and two pumps. These actuators
1519 regulate the level within the tank by controlling the inflow and out-
1520 flow of the liquid. These sensors are identified as MV101 (valve), P101
1521 and P101 (pumps) in the PLC registers.

1522 Despite its moderate complexity, this testbed provides an ideal plat-
1523 form for presenting straightforward and concise examples of the frame-
1524 work's behavior. It enables us to effectively demonstrate the potential
1525 of the framework and facilitate a deeper understanding of its underlying
1526 methodology.

1527 4.2.2 Phase 0: Network Analysis

1528 The objective of the network analysis presented in this section is to offer
1529 users valuable information regarding the communication process within
1530 an industrial control system and a broader perspective on network com-
1531 munications, delving into previously unexplored aspects of the system.
1532 These additional dimensions provide a deeper understanding of the sys-
1533 tem's behavior and characteristics. This analysis aims to provide users
1534 with an overview of the communication between PLCs at level 1 (see Fig-
1535 ure 2.1), as well as the communication between PLCs and devices at higher
1536 levels such as HMIs and Historian servers (see Section 2.1 for ICSs archi-
1537 tecture). This also allows us to have a better understanding of the system
1538 architecture and network topology. The analysis focuses on industrial pro-
1539 tocols used and the information exchanged.

1540 By reconstructing the network communication structure using data ob-
1541 tained from the network traffic sniffing process, users can gain a compre-
1542 hensive understanding of the behavior of the underlying industrial sys-

1543 tem. This knowledge can then be utilized to **plan a strategy** for analyzing
1544 the physical processes within the system.

1545 **4.2.2.1 Extracting Data from PCAP Files**

1546 The initial step involves extracting the desired information from the
1547 PCAP files that contain the captured network traffic. This includes details
1548 such as the source IP address, destination IP address, protocol used, and
1549 the type of request made (e.g., Read/Write, Request/Response). The ex-
1550 tracted data is then converted into a more convenient CSV format. This
1551 extracted data serves as the foundation for the subsequent phase.

1552 In the latter part of this phase, the extracted data is utilized to generate
1553 the network schema. The network schema provides a visual representa-
1554 tion of the connections and relationships within the network, showcasing
1555 the communication patterns between different components. This schema
1556 helps in understanding the overall structure and behavior of the industrial
1557 control system.

1558 To accomplish the extraction of data from the PCAP files, a Python
1559 script called `export_pcap_data.py` is employed. This script, originally de-
1560 signed for the business process phase, is located in the directory
1561 `$(project_dir)/network-analysis` and accepts the following options as
1562 command-line arguments:

- 1563 • **-f or --filename:** allows the user to specify a single PCAP file to be
1564 passed as input to the script. The user can provide the complete file
1565 path of the PCAP file as an argument;
- 1566 • **-m or --mergefiles:** enables the merging of multiple PCAP files. In
1567 this scenario, the files should be located within the directory speci-
1568 fied by the `pcap_dir` directive in the `config.ini` configuration file and
1569 the user does not have to provide the path to each PCAP file;
- 1570 • **-d or --mergedir:** allows for specifying the directory that contains the
1571 PCAP files to be automatically imported into the script and merged.

1572 This ensures that all the PCAP files within the specified directory will
1573 be processed by the script without the need for manual selection or
1574 input.

- 1575 • **-s or --singledir:** operates differently from the previous option men-
1576 tioned. This option enables the extraction of data from each individ-
1577 ual PCAP file within the specified directory. The extracted data is
1578 then saved in separate CSV datasets, which are stored in the direc-
1579 tory specified by the `split_dir` directive in the `config.ini` file. This
1580 functionality proves useful when dealing with exceptionally large
1581 PCAP files, where merging them together for export might consume
1582 significant time and resources. By utilizing this option, the extrac-
1583 tion procedure becomes lighter and more manageable. The extracted
1584 data in separate CSV files can be utilized in the later stages of the
1585 Network Analysis process;
- 1586 • **-t or --timerange:** this functionality enables users to specify a specific
1587 time period within the PCAP files from which they wish to extract
1588 relevant information.

1589 Unless the `-s` option is explicitly specified, the results of data extraction
1590 and export will be saved to a single CSV dataset within the
1591 `$(project_dir)/network-analysis/data` directory. The default file name
1592 for this output file is determined by the `pcap_export_output` directive spec-
1593 ified in the `config.ini` file. In addition, by utilizing the `protocols` and
1594 `ws_<protocol>_field` directives, user can configure the network protocols
1595 to be searched within the PCAP files. Furthermore, user can specify the
1596 relevant Tshark/Wireshark fields to extract for the specified protocols set
1597 in the `protocols` directive.

1598 After obtaining the extracted data, it is possible to proceed with the
1599 second part of the network analysis.

1600 4.2.2.2 Network Information

1601 During this stage, the exported CSV data is processed to derive valua-
1602 able information regarding network communications and the structure of
1603 the network itself. The objective is to identify and establish relationships
1604 between IP addresses present on the network, thereby determining the
1605 sources and destinations of communications. Furthermore, the analysis
1606 detects the protocols used for each communication and quantifies the var-
1607 ious types of requests made.

1608 This information is then transformed into a **graph representation of**
1609 **the network** (or subnetwork, if specified). In this graph, devices are repre-
1610 sented as nodes labeled with their IP addresses, while edges represent the
1611 incoming and outgoing communications of these devices, along with the
1612 corresponding information.

1613 To ensure comprehensibility, the analysis also provides users with **tex-**
1614 **tual information** containing the same details as the graph representation.
1615 This text-based information serves as an alternative for cases where the
1616 graph may become complex to interpret, particularly when numerous edges
1617 connect nodes and result in a high volume of network requests.
1618 This textual information is saved to another CSV file, enabling offline ref-
1619 erence or potential future utilization. By having this file available, users
1620 can access the network analysis results in a structured format for further
1621 analysis or documentation purposes.

1622 The Python script `networkAnalysis.py` in the `$(project_dir)/network-analysis`
1623 directory manages this phase of the analysis. The script can be executed
1624 with the following parameters:

- 1625 • **-f or --filename:** used to specify the CSV dataset containing the net-
1626 work data exported in the previous step. The dataset should be lo-
1627 cated in the directory `$(project_dir)/network-analysis/data`;
- 1628 • **-D or --directory:** used to specify the directory that contains the CSV
1629 datasets obtained using the `-s` option of the Python script `export_pcap_data.py`.

1630 By passing this parameter, the script will automatically merge the
1631 datasets and proceed with the analysis of the data contained within
1632 them;

- 1633 • **-s or --srcaddr:** allows for specifying the source IP address for which
1634 you wish to display the incoming and outgoing communications. By
1635 providing the source IP address as an argument, the script will focus
1636 on showcasing the communications associated with that particular
1637 IP address;
- 1638 • **-d or --dstaddr:** enables the user to specify the destination IP address
1639 for which you want to display the incoming and outgoing communi-
1640 cations. By providing the destination IP address as an argument, the
1641 script will concentrate on presenting the communications associated
1642 with that specific IP address.

1643 The parameters related to IP addresses, including source and destina-
1644 tion, are optional. It is possible to specify either one of them individually.
1645 For instance, if the user specifies only the source IP address, the script will
1646 display the network nodes with which it communicates on the outgoing
1647 side, along with the corresponding generated traffic. Similarly, if only the
1648 destination IP address is specified, the script will showcase the network
1649 nodes communicating with it on the incoming side, along with the rele-
1650 vant traffic data.

1651 During the analysis, the script identifies and displays the IP addresses
1652 present in the network as output for the user's reference. This allows the
1653 user to select specific IP addresses from the command line for a more fo-
1654 cused analysis, such as choosing a subnet of interest. Additionally, the
1655 script detects and tracks distinct communications between pairs of PLCs,
1656 keeping a record of the number of these communications.

1657 The result of the analysis is a graph representation of the network (or
1658 subnetwork) to be analyzed. An example of such a graph can be seen in
1659 Figure 4.1.

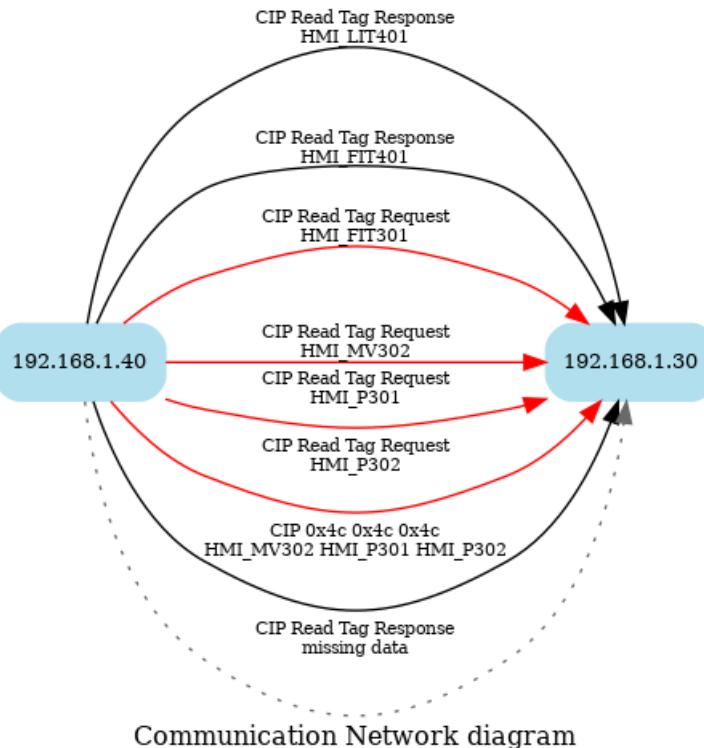


Figure 4.1: Network Communications between 192.168.1.40 (Source) and 192.168.1.30 (Destination)

1660 The graph illustrates the network communications between the source
 1661 IP address 192.168.1.40 and the destination IP address 192.168.1.30. Each
 1662 arrow represents a communication between these IP addresses, showcas-
 1663 ing the flow and direction of the interactions. The red arrows indicate re-
 1664 quests initiated by the source IP address towards the destination IP, while
 1665 the black arrows represent responses sent by the source IP in response
 1666 to previous requests made by the destination IP. The gray dotted arrows
 1667 represent responses for which the corresponding request is missing or un-
 1668 available for some reason. Overall, the graph distinguishes the different
 1669 types of communications and provides insights into the request-response
 1670 dynamics between the source and destination IP addresses. The graph is
 1671 automatically generated and saved within the
 1672 `$(project_dir)/network-analysis/data` directory.

1673 In terms of the textual output, we can observe how the same data is
 1674 represented. The communications exchanged between the two PLCs are
 1675 displayed more prominently, allowing for a clearer understanding. Unlike
 1676 the graph, the textual representation includes a column on the right-hand
 1677 side, indicating the number of communications for each type of request.
 1678 This provides a more distinct perspective on the network behavior within
 1679 an industrial control system that utilizes, in this case, the CIP protocol for
 1680 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

1681 As previously mentioned, the textual data is stored in a CSV dataset
 1682 located in the \$(project_dir)/network-analysis/data directory.

1683 4.2.3 Phase 1: Data Pre-processing

1684 *Data Pre-processing phase* is probably the most delicate and significant
 1685 one: depending on how large the industrial system to be analyzed is, the
 1686 data collected, and how it is enriched using the additional attributes, the
 1687 subsequent system analysis will provide more or less accurate outcomes.

1688 The Ceccato et al.'s tool has several limitations, particularly at this
 1689 stage. It does not allow for the isolation of a subsystem, either in terms
 1690 of time or the number of PLCs to be analyzed. The system is considered as
 1691 a whole without the ability to focus on specific subsystems. Additionally,
 1692 many of the additional attributes had to be manually added, and for the
 1693 ones entered automatically, there is no way to specify the register type to
 1694 associate them with.

1695 The combination of these limitations, along with the presence of hard-
1696 coded references to attributes and registers in the tool’s code, makes the
1697 analysis of the system more challenging. Furthermore, it compromises the
1698 accuracy and reliability of the obtained results in terms of both quantity
1699 and quality.

1700 In the proposed framework, these issues have been addressed by in-
1701 corporating new features.

1702 *Firstly*, the framework allows for the **selection of a subsystem from**
1703 **the command line** based on both temporal criteria and the specific PLCs
1704 to be included. This enables more focused and targeted analysis.

1705 *Secondly*, we have **revamped the process of enriching** the resulting
1706 dataset by eliminating manual entry of additional attributes. Instead, users
1707 now have the flexibility to determine the type of additional attribute to as-
1708 sociate with a specific register.

1709 *Thirdly*, after the pre-processing stage, a **preliminary analysis** can be
1710 conducted on the resulting dataset. This analysis aims to identify the reg-
1711 isters that are associated with actuators, measurements, and hardcoded
1712 setpoints or constants. It provides insights into the dataset and helps in
1713 refining the enrichment step. The parameters for this analysis can be con-
1714 figured in the *config.ini* file, allowing for customization and fine-tuning of
1715 the process.

1716 In the upcoming sections, we will delve into a more comprehensive
1717 examination of the achievements made in this phase.

1718 4.2.3.1 Subsystem Selection

1719 In Ceccato et al.’s tool, the datasets for each individual PLC in CSV for-
1720 mat were required to be placed in a specific directory that was hardcoded
1721 in the script. The script would then merge and enrich these datasets to
1722 generate a single output dataset representing the complete process trace of
1723 the industrial system. However, the script did not provide options to se-
1724 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:

```

1742 [PATHS]
1743 root_dir = /home/marcuzzo/UniVr/Tesi
1744 project_dir = %(root_dir)s/PLC-RE
1745 net_csv_path = %(root_dir)s/datasets_SWaT/2015/Network_CSV
1746
1747 [PREPROC]
1748 raw_dataset_directory = datasets_SWaT/2015 # Directory
1749     ↪ containing datasets
1750 dataset_file = PLC_SWaT_Dataset.csv # Default output
1751     ↪ dataset
1752 granularity = 10 # slope granularity
1753 number_of_rows = 20000 # Seconds to consider
1754 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

1758 through the command line will override the default settings specified in
1759 the *config.ini* file. These options are:

- 1760 • **-s or --skiprows:** initial transient period (expressed in seconds) to
1761 be skipped. This option is useful in case the system has an initial
1762 transient or the analyzer wishes to start the analysis from a specific
1763 point in the dataset;
- 1764 • **-n or --nrows:** time interval under analysis, expressed in terms of the
1765 number of rows in the dataset.
1766 This option makes a **selection** on the data of the dataset;
- 1767 • **-p or --plcs:** PLCs to be merged and enriched. The user can specify
1768 the desired PLCs by indicating the CSV file names of the associated
1769 datasets with no limitations on number.
1770 This option makes a **projection** on the data of the dataset.
- 1771 • **-d or --directory:** performs the merge and enrichment of all CSV files
1772 contained in the directory specified by user, overriding the default
1773 setting in *config.ini*. It is in fact the old functionality of the Ceccato
1774 et al.'s tool, maintained here to give the user more flexibility and
1775 convenience in case he wants to perform the analysis on the whole
1776 system. This is also the default behavior in case the -p option is not
1777 specified.
- 1778 • **-o or --output:** specifies the name of the file in which the obtained
1779 dataset will be saved. It must necessarily be a file in CSV format.
- 1780 • **-g or --granularity:** specifies a granularity (expressed in seconds)
1781 that will be used to calculate the measurement slope during the dataset
1782 enrichment phase. We will discuss this later in Section 4.2.3.2.

1783 **4.2.3.2 Dataset Enrichment**

1784 After a step in which a function is applied to each PLC-related dataset
 1785 to eliminate its unused registers within the system¹, the **dataset enrichment operation** is performed.
 1786

1787 This operation brings about several distinctions compared to the previous version. Firstly, it is performed on each individual dataset instead of the resulting dataset. Additionally, there are a greater number of attributes, which are automatically calculated and added to the dataset by the `mergeDatasets.py` script. Importantly, the configuration file `config.ini` under the [DATASET] section allows users to determine which registers should be assigned to these attributes.

1794 In Listing 4.3 we can see the list of additional attributes and how they should be associated with the registers of the dataset:
 1795

```
1796 [DATASET]
1797 timestamp_col = Timestamp
1798 max_prefix = max_
1799 min_prefix = min_
1800 max_min_cols_list = lit|ait|dpit
1801 prev_cols_prefix = prev_
1802 prev_cols_list = mv[0-9]{3}|p[0-9]{3}
1803 trend_cols_prefix = trend_
1804 trend_cols_list = lit
1805 trend_period = 150
1806 slope_cols_prefix = slope_
1807 slope_cols_list = lit
```

Listing 4.3: config.ini parameters for dataset enriching

1808 In the following, we report a brief explanation of the parameters just seen:

1809 **timestamp_col** denotes the name of the column in the dataset that holds
 1810 the timestamp data. This parameter is significant not only in the 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.

1812 al.'s work, this parameter was hardcoded and lacked configurability,
1813 leading to errors if the analyzed system changed.

1814 **max_prefix, min_prefix, max_min_cols_list** refer to any relative maximum
1815 or minimum values (*relative setpoints*) of one or more measures and
1816 that can be found and inserted as new columns within the dataset.
1817 The first two parameters indicate the prefix to be used in the column
1818 names affected by this additional attribute, while the third specifies
1819 of which type of registers we want to know the maximum and/or
1820 minimum value reached (several options can be specified using the
1821 logical operator | - or).

1822 If, for example, we want to know the maximum value of the registers
1823 associated with the tanks, indicated in the iTrust SWaT system by the
1824 prefix LIT, we only need to specify the necessary parameter in the
1825 *config.ini* file, so `max_min_cols_list = lit`.

1826 The result will be to have in the dataset thus enriched a new column
1827 named `max_LIT101`.

1828 **prev_cols_prefix, prev_cols_list** refer to the values at the previous time
1829 instant of the registers specified in `prev_cols_list`. It is possible to
1830 specify registers using *regex*, as in the example shown. It may be use-
1831 ful in some cases to have this value available to check, for example,
1832 when a change of state of a single given actuator occurs. The behav-
1833 ior of these parameters is the same as described in the point above.

1834 **slope_cols_prefix, slope_cols_list** are related to the calculation of the
1835 slope of a specific register that contains numeric values (usually a
1836 measure), that is, its trend. Slope calculation makes little sense on
1837 booleans. The slope can be **ascending** (if its value is greater than
1838 zero), **descending** (if less than zero) or **stable** (if approximately equal
1839 to zero). We will delve into the details of slope calculation in the fol-
1840 lowing paragraph, as it pertains to the attributes `trend_cols_prefix`,
1841 `trend_cols_list`, and `trend_period`.

1842 Initially, the parameters for registers to be associated with each addi-
1843 tional attribute may be left blank, as we may not have prior knowledge
1844 about the system and are unsure about which registers correspond to ac-
1845 tuators, measurements, or other attributes. This information can be ob-
1846 tained from the preliminary analysis that follows the merging of datasets.
1847 The analysis, performed based on user's choice, provides indications on
1848 potential sensors, actuators, and other relevant information. These indica-
1849 tions help the user set the desired values in the *config.ini* file and refine the
1850 enrichment process by re-launching the *mergeDatasets.py* script.

1851 **Slope Calculation** The *slope* is an attribute that represents the **trend** of
1852 the measurement being considered. It is particularly useful, in our con-
1853 text, during the inference and invariant analysis phase to gather informa-
1854 tion about the trend under specific conditions. The slope can generally be
1855 classified as **increasing** (*slope* > 0), **decreasing** (*slope* < 0), or **stable** (*slope*
1856 = 0).

1857 Normally, the slope is calculated through a simple mathematical for-
1858 mula: given an interval *a*, *b* relative to the measurement *l*, the slope is
1859 given by the difference of these two values divided by the amount of time
1860 *t* that the measurement takes to reach *b* from *a*:

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

1861 In the proposed framework, similar to the Ceccato et al.'s tool, this time
1862 interval (the granularity) can be adjusted to be either long or short.

1863 The choice of granularity depends on the desired accuracy of the slope
1864 calculation. A lower granularity will provide a slope that closely reflects
1865 the actual measurement trend, while a higher granularity will result in
1866 flatter slope data. Each time interval within which the measurement is
1867 divided corresponds to a slope value. These slopes are calculated and
1868 added as additional attributes in the dataset. Later on, these slope values
1869 are used to determine the trend of the measurement in specific situations

1870 or conditions.

1871 Calculating the slope directly from the raw measurement data can be
1872 a suitable approach for systems where the measurements are *not* heavily
1873 influenced by **perturbations**. Perturbations, such as liquid oscillations in
1874 a tank during filling and emptying phases, can lead to fluctuating read-
1875 ings of the level. In such cases, maintaining a low granularity can provide
1876 a more accurate calculation of the overall trend that closely aligns with
1877 the actual measurement trend. This situation occurs, for instance, in the
1878 testbed utilized by Ceccato et al.

1879 However, if perturbations significantly affect the measurement read-
1880 ings, calculating the slope on individual time intervals may result in an
1881 inaccurate trend definition, irrespective of the chosen granularity. In such
1882 cases, the fluctuating nature of the measurements due to perturbations can
1883 introduce errors in the slope calculation, making it less reliable as an indi-
1884 cator of the actual trend.

1885 Figure 4.3 demonstrates this assertion: the measurement, in blue, refers
1886 to the LIT101 tank of our testbed; in red, the slope calculation related to
1887 the measurement with three different granularities: 30 (Figure 4.3a), 60
1888 (Figure 4.3b) and 120 seconds (Figure 4.3c).

1889 It is noticeable that as the granularity increases, the slope values flatten.
1890 Moreover, in the time interval between seconds 1800 and 4200, the level of
1891 LIT101 exhibits a predominantly increasing trend, yet the calculated slope
1892 values fluctuate between positive and negative. Consequently, during the
1893 invariant analysis, the overall increasing trend may not be detected, re-
1894 sulting in a loss of information.

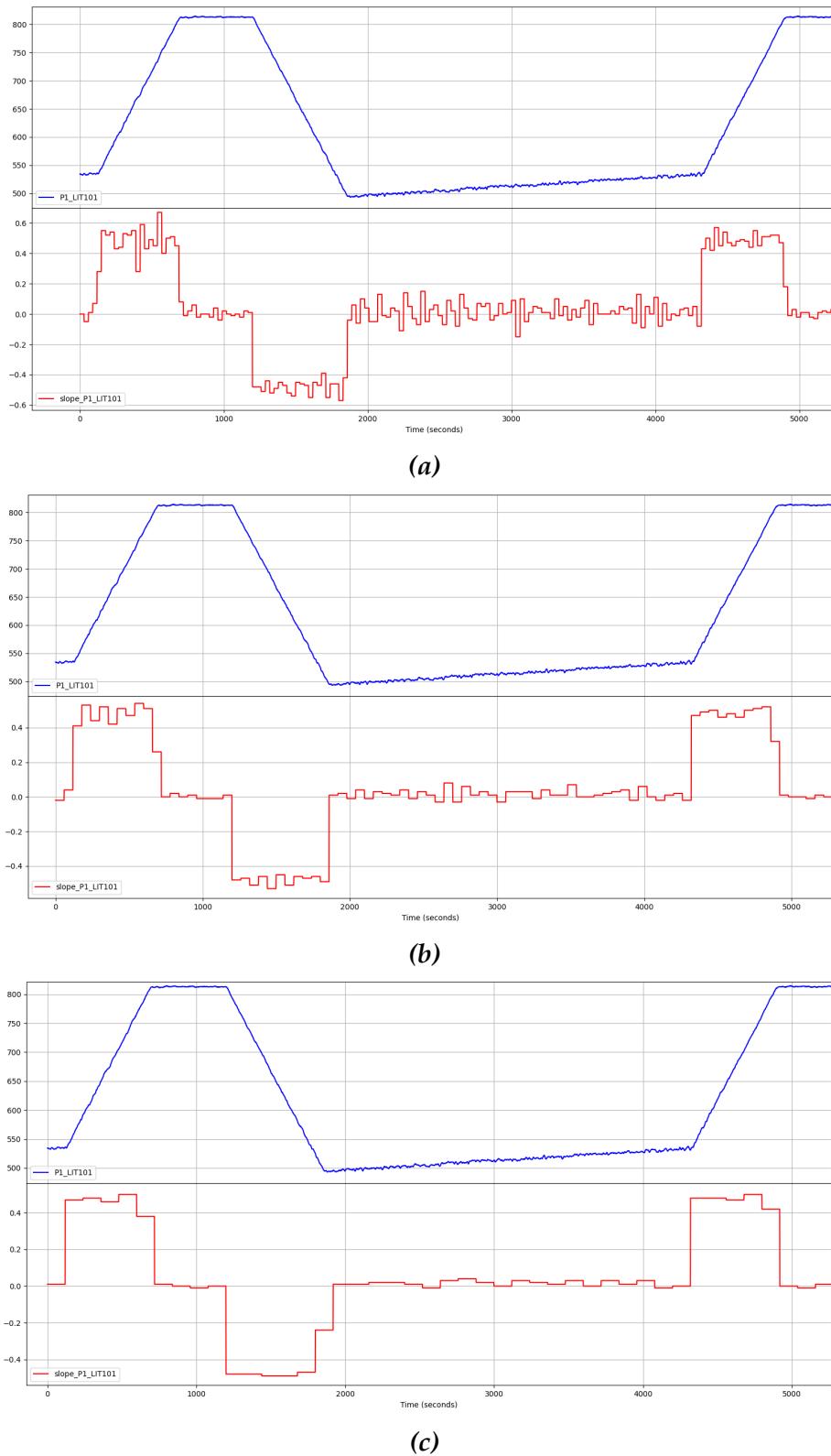


Figure 4.3: Slope comparison with granularity 30 (a), 60 (b) and 120 seconds (c)

1895 Ceccato et al.'s tool did not take into account the possibility of having
1896 strongly perturbed data, which presented a challenge that we needed to
1897 address in the development of the proposed framework.

1898 The solution to this problem involves applying techniques to reduce
1899 the "noise" in the data, aiming to achieve a more linear trend in the mea-
1900 surement curve. By minimizing the effects of perturbations, we can calcu-
1901 late slopes more accurately.

1902 There are various methods available for smoothing out noise in the
1903 data. In our framework, we focused on two commonly used approaches
1904 found in the literature: **polynomial regression** and **seasonal decomposi-**
1905 **tion**. In addition to these two methods, we also explored the use of a **line**
1906 **simplification algorithm**.

1907 *Polynomial regression* [58] is a technique that allows us to create a filter
1908 to reduce the impact of noise on the data. By fitting a polynomial function
1909 to the measurements, we can obtain a smoother curve that captures the
1910 underlying trend while minimizing the effects of perturbations.

1911 *Seasonal decomposition* [59], specifically the part related to trending, is
1912 another method we explored. It involves decomposing the time series into
1913 different components, such as trend, seasonality, and residual. By isolat-
1914 ing the trend component, we can obtain a cleaner representation of the
1915 underlying pattern in the data.

1916 *Line simplification algorithms* [60] aim to reduce the complexity of a poly-
1917 line or curve by approximating it with a simplified version composed of
1918 fewer points. By selectively removing redundant or less significant points,
1919 line simplification algorithms help reduce storage space and computa-
1920 tional requirements while preserving the overall shape and characteristics
1921 of the original line.

1922 Regarding polynomial regression, we evaluated the use of the **Savitzky-**
1923 **Golay filter** [61] as a smoothing technique. For seasonal decomposition,
1924 we explored the **Seasonal-Trend decomposition using LOESS** (STL) method
1925 [62]. For the line simplification algorithm, we specifically considered the

¹⁹²⁶ **Ramer-Douglas-Peucker (RDP) algorithm [63].**

¹⁹²⁷ Figure 4.4 shows a quick graphical comparison of these techniques
¹⁹²⁸ compared with the original data. The solution adopted is the *STL decomposition*
¹⁹²⁹ method, which effectively reduces noise compared to the Savitzky-
¹⁹³⁰ Golay filter. However, it should be noted that this method may introduce
¹⁹³¹ some delay in certain parts of the data, as is typically observed in similar
¹⁹³² algorithms. Despite its apparent effectiveness, the RDP algorithm fails to
¹⁹³³ accurately approximate sections where the measurement level remains rel-
¹⁹³⁴ atively stable. Consequently, it yields incorrect slope estimations, causing
¹⁹³⁵ a loss of valuable information about the system.

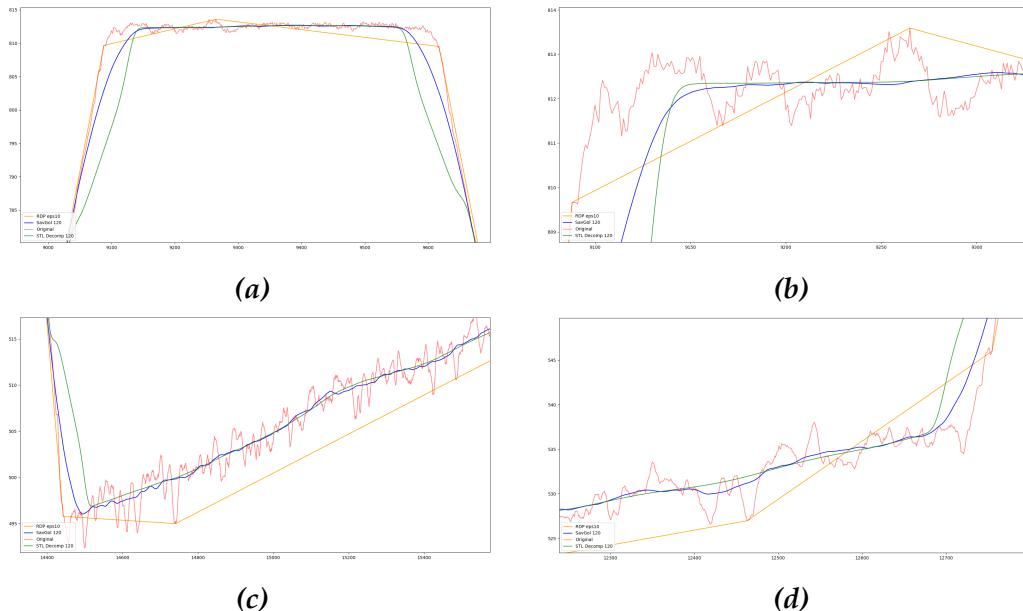


Figure 4.4: Savitzky-Golay filter (blue line), STL decomposition (green) and RDP algorithm (orange) comparison

¹⁹³⁶ By applying the STL decomposition, we observe a notable enhance-
¹⁹³⁷ ment in slope calculation even when using a low granularity. Figure 4.5
¹⁹³⁸ demonstrates that, with the same granularity as shown in Figure 4.3a, the
¹⁹³⁹ slope values, albeit exhibiting fluctuations, consistently align with the un-
¹⁹⁴⁰ derlying trend of the data curve. The introduced lag resulting from the
¹⁹⁴¹ decomposition's periodicity is responsible for the observed delay.

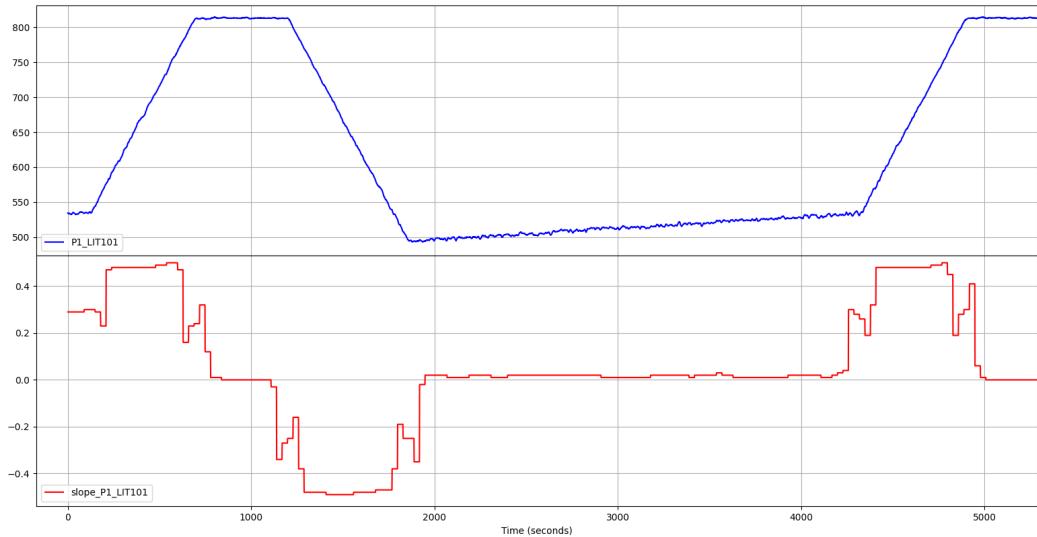


Figure 4.5: Slope after the application of the STL decomposition

1942 The periodicity, which defines the sampling time window for decom-
 1943 position and the level of noise smoothing, can be configured using the
 1944 trend_period directive in the *config.ini* file.
 1945 During the slope calculation, the analysis will be performed on the data
 1946 from the additional measurement trend attributes specified in the trend_cols_list
 1947 directive of the configuration file, rather than on the original unfiltered
 1948 data.

1949 To ensure proper interpretation by Daikon, the decimal values repre-
 1950 senting the calculated slopes are converted into **three numerical values**:
 1951 -1, 0, and 1. These values correspond to *decreasing* (if the slope is less than
 1952 zero), *stable* (if it is equal to zero), and *increasing* (if it is greater than zero)
 1953 trends, respectively. Figure 4.6 displays the modified slopes along with
 1954 the curve obtained from the STL decomposition:

1955 4.2.3.3 Datasets Merging

1956 During this step, the datasets of the individual PLCs are merged, re-
 1957 sulting in two separate datasets. The first dataset is enriched with addi-
 1958 tional attributes but excludes the timestamp column. This dataset is in-

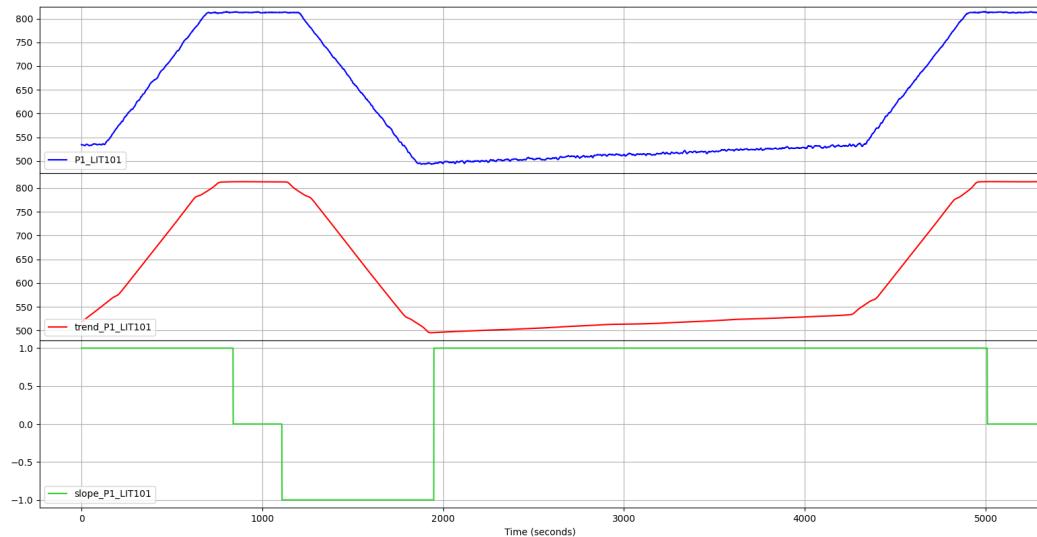


Figure 4.6: The new slope representation (green line) and the smoothed measurement data obtained with the STL decomposition (red)

1959 tended for inference and invariant analysis. The second dataset does *not*
 1960 contain any additional data and is specifically used in the process mining
 1961 phase. This dataset contains the timestamp.

1962 By default, the enriched dataset will be saved in CSV format in the
 1963 `$(project-dir)/daikon/Daikon_Invariants` directory. The other dataset,
 1964 without additional data, will be saved in the `$(project-dir)/process-mining/data`
 1965 directory. It's worth noting that both paths can be configured in the
 1966 `config.ini` file. The dataset name can be specified in the `config.ini` file or
 1967 through the `-o` command-line option. When generating the dataset for
 1968 process mining, the script will automatically add a `_TS` suffix to the file-
 1969 name to indicate that it includes the timestamp. This flexibility allows the
 1970 user to provide a different filename for each output, preventing overwrit-
 1971 ing of previous datasets. It enables the user to save the execution trace of
 1972 the selected subsystem separately and utilize them in subsequent analysis
 1973 phases.

1974 4.2.3.4 Preliminary Analysis of the Obtained Subsystem

1975 After merging the datasets, the user has the option to perform an **optional analysis** of the resulting dataset to extract preliminary data. This
 1976 analysis aims to gather basic information about the (sub)system and po-
 1977 tentially refine the enrichment process. If the user chooses to proceed with
 1978 the analysis, the `mergeDatasets.py` script invokes another Python script lo-
 1979 cated in the `$(project-dir)/pre-processing` directory called `system_info.py`.
 1980

1981 Relying on an analysis based on a combination of Daikon and Pandas
 1982 this script performs a quick analysis of the dataset allowing to **estimate**, al-
 1983beit approximately, the **type of registers** (sensors, actuators, ...), also iden-
 1984tifying possible maximum and minimum values of measurements and
 1985 hardcoded setpoints. Furthermore, leveraging the use of the additional
 1986 attribute `prev_`, the `system_info.py` script is capable of deriving measure-
 1987 ment values corresponding to state changes of individual actuators. This
 1988 allows for the identification of specific measurements associated with the
 1989 activation or deactivation of certain actuators within the system.

1990 As the last information we have duration of actuator states for each
 1991 cycle of the system: this information can be useful for making assumptions
 1992 and conjectures about the behavior of an actuator in a specific state or, by
 1993 observing the duration values of each cycle, highlighting anomalies in the
 1994 system.

1995 Listing 4.4 shows an example of the output this brief analysis related to
 1996 our testbed (for brevity, only one measurement is reported in the analysis
 1997 of actuator state changes):

```
1998 Do you want to perform a brief analysis of the dataset? [y
1999 ↵ /n]: y
2000
2001 Actuators:
2002 MV101 [0.0, 1.0, 2.0]
2003 P101 [1.0, 2.0]
2004
2005 Sensors:
2006 FIT101 {'max_lvl': 2.7, 'min_lvl': 0.0}
```

```
2007     LIT101 {'max_lvl': 815.1, 'min_lvl': 489.6}
2008
2009     Hardcoded setpoints or spare actuators:
2010     P102 [1.0]
2011
2012     Actuator state changes:
2013     LIT101      MV101      prev_MV101
2014     800.7170    0          2
2015     499.0203    0          1
2016     800.5992    0          2
2017     498.9026    0          1
2018     800.7170    0          2
2019     499.1381    0          1
2020     801.3058    0          2
2021     498.4315    0          1
2022     801.4628    0          2
2023     498.1567    0          1
2024
2025     LIT101      MV101      prev_MV101
2026     805.0741    1          0
2027     805.7414    1          0
2028     805.7806    1          0
2029     805.1133    1          0
2030     804.4068    1          0
2031
2032     LIT101      MV101      prev_MV101
2033     495.4483    2          0
2034     497.9998    2          0
2035     495.9586    2          0
2036     495.8016    2          0
2037     494.5847    2          0
2038
2039     LIT101      P101      prev_P101
2040     536.0356    1          2
2041     533.3272    1          2
2042     542.1591    1          2
2043     534.8581    1          2
2044     540.5890    1          2
2045
```

```

2046      LIT101        P101        prev_P101
2047      813.0031       2           1
2048      813.0031       2           1
2049      811.8256       2           1
2050      812.7283       2           1
2051      813.3171       2           1
2052
2053      Actuator state durations:
2054      MV101 == 0.0
2055      9   9   10   9   9   10   9   9   10   9
2056
2057      MV101 == 1.0
2058      1174  1168  1182  1160  1172
2059
2060      MV101 == 2.0
2061      669   3019  3012  3000  2981
2062
2063      P101 == 1.0
2064      1073  1074  1061  1056  1056
2065
2066      P101 == 2.0
2067      118   3133  3143  3125  3108

```

Listing 4.4: Example of preliminary system analysis

2068 The information obtained here can be used, as mentioned above, to
 2069 refine the enrichment of the dataset by setting directives in the [DATASET]
 2070 section of the *config.ini* file, should this be empty or only partially set, or to
 2071 make the first conjectures about the system, as we have just seen.

2072 The *system_info.py* file can also run in standalone mode if needed:
 2073 it takes as command-line arguments the dataset to be analyzed, a list of
 2074 actuators, and a list of sensors. For analysis related to state changes, the
 2075 dataset must mandatorily be of the enriched type.

2076 4.2.4 Phase 2: Graphs and Statistical Analysis

2077 The introduction of the new *graph analysis* is motivavted by from the re-
 2078 quirement to provide users with a comprehensive overview of the (sub)system
 2079 derived from the preceding pre-processing phase. The objective is to facil-
 2080 itate the identification of register types, enhance the understanding of re-
 2081 lationships, and effectively grasp the dynamics among registers controlled
 2082 by one or more PLCs. This analysis component serves to validate initial
 2083 conjectures made during the preliminary analysis described in the previ-
 2084 ous section or generate new insights with the aid of visual chart represen-
 2085 tations. It enables users to gain a deeper understanding of the system by
 2086 leveraging the support of visual charts.

2087 In Ceccato et al.'s framework, as mentioned in Section 3.2.7, it was only
 2088 possible to view the chart of one register at a time. While this allowed for
 2089 the identification or hypothesis of the register type, it made it challenging
 2090 to establish relationships with other components of the system and derive
 2091 conjectures about their behavior. To address this limitation, there was a
 2092 need for a new tool that could provide more information in a more acces-
 2093 sible manner.

2094 Initially, we considered adopting an approach similar to Figure 3.5,
 2095 where all the graphs are displayed within a single plot. However, we soon
 2096 realized that this solution was not feasible and could not be adopted. Fig-
 2097 ure 4.7 helps to understand why.

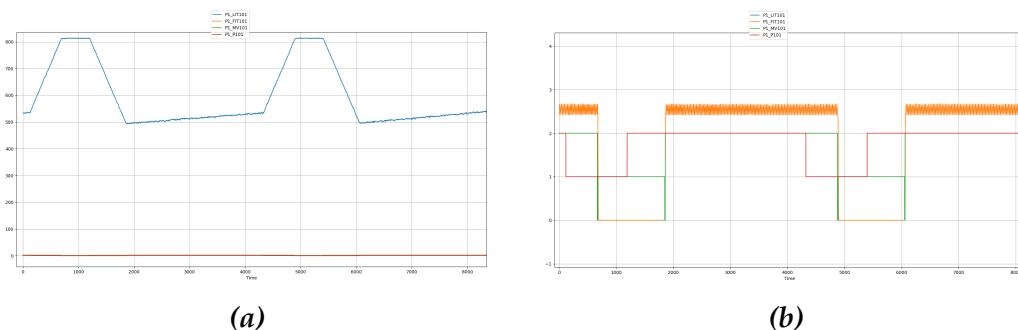
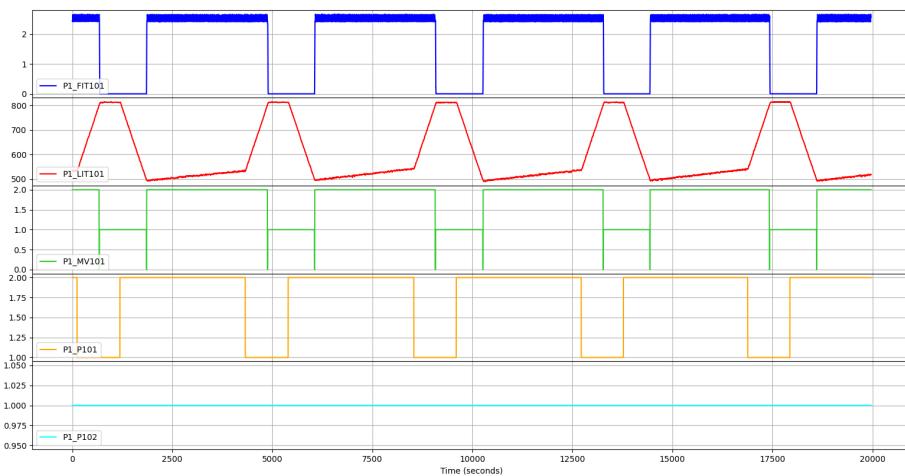


Figure 4.7: Plotting registers on the same y-axis

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.

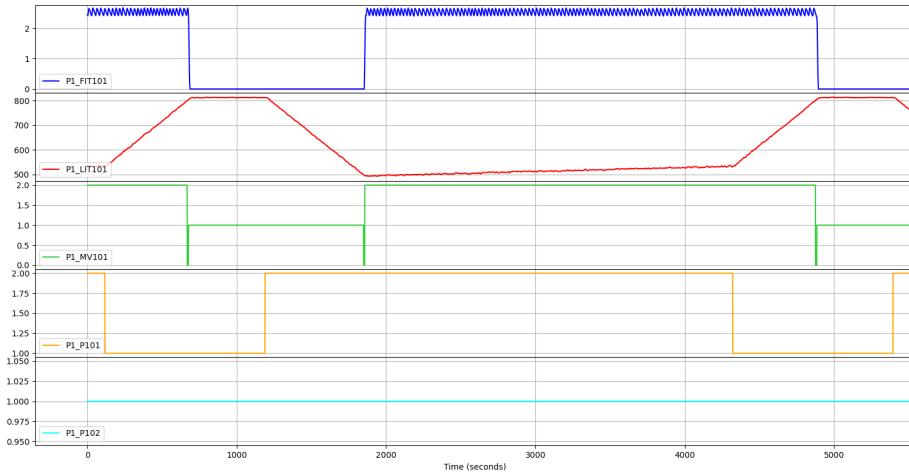
The solution to this issue is simple and effective: the use of **subplots**. Basically, each register corresponds to a subplot of the graph that shares the time axis (the x-axis) with the other subplots, but keeps the y-axis of the values of each register independent. This maintains the readability and 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.)

2115 As the reader has probably already noticed, the majority of the graphs
 2116 presented in the previous sections and chapters were generated using the
 2117 new graph analysis script, specifically the `runChartsSubPlots.py` script.
 2118 This Python script is located in the `$(project-dir)/statistical-graphs`
 2119 directory and utilizes the `matplotlib` libraries to generate the graph plots,
 2120 similar to the Ceccato et al.'s tool.

2121 The script accepts the following command-line parameters:

- 2122 • **-f or --filename:** specifies the CVS format dataset to read data from.
 2123 The dataset must be within the directory containing the enriched
 2124 datasets for the invariant analysis phase. If subsequent parameters
 2125 are not specified, the script will display all registers in the dataset,
 2126 excluding any additional attributes;
- 2127 • **-r or --registers:** specifies one or more specific registers to be dis-
 2128 played;
- 2129 • **-a or --addregisters:** adds one or more registers to the default visual-
 2130 ization. This option is useful in case an additional attribute such as
 2131 slope is to be analyzed;

- 2132 • **-e or --excluderegisters:** excludes one or more specific registers from
2133 the default visualization. This option is useful to avoid displaying
2134 hardcoded setpoints or spare registers.

2135 This script, like the previous ones, is designed to provide the maximum
2136 flexibility and ease of use for the user, combined with greater power and
2137 effectiveness in deriving useful information about the analyzed system.

2138 **Statistical Analysis** After careful consideration, we made the decision
2139 not to include the statistical analysis aspect of the Ceccato et al.'s tool in
2140 the framework. We found that there was no practical use for it. Instead, we
2141 integrated the relevant statistical information into the preliminary analysis
2142 conducted after the pre-processing phase. Additionally, we deemed the
2143 histogram to have limited utility and considered it outdated in comparison
2144 to the new graph analysis approach we implemented.

2145 Although we decided not to include the statistical analysis aspect in
2146 the framework, the Python script `histPlots_Stats.py` from the original
2147 tool remains in the directory. This script is essentially unchanged from the
2148 version developed by Ceccato et al. and can be used in the future if the
2149 need arises.

2150 4.2.5 Phase 3: Invariant Inference and Analysis

2151 The phase of invariant inference and analysis has undergone a redesign
2152 and improvement to offer the user a more comprehensive and easier ap-
2153 proach to identify invariants. This has been achieved through the applica-
2154 tion of new criteria to analyze and reorganize the Daikon analysis results.
2155 The outcome of this is a more compact presentation of information that
2156 highlights the possible relationships among invariants.

2157 The new design not only enables the identification of undiscovered as-
2158 pects of the system behavior but also confirms the hypotheses made dur-
2159 ing the earlier stages of analysis. This step is now semi-automated, unlike
2160 before, and allows for the analysis of invariants on individual actuator
2161 states and their combinations.

2162 **4.2.5.1 Revised Daikon Output**

2163 To streamline the process of identifying invariants quickly and effi-
2164 ciently, it is necessary to revise the output generated by standard Daikon
2165 analysis. The goal is to create a more compact and readable format for the
2166 output.

2167 The current Daikon results basically consist of three sections, referred
2168 as *program point sections* [64]:

- 2169 1. the first section containing **general invariants**, i.e., valid regardless
2170 of whether a condition is specified for the analysis;
- 2171 2. the second section containing **conditional invariants**, obtained by
2172 specifying conditions for the analysis in the *.spininfo* file.
- 2173 3. a third section containing the invariants that are obtained from the
2174 **negation of the condition** potentially specified in the *.spininfo* file.

2175 In each section only a single invariant per row is shown, without relat-
2176 ing it in any way to the others: this makes it difficult to identify significant
2177 invariants and any invariant chain that might provide much more infor-
2178 mation about the behavior of the system than the single invariant.

2179 A brief example of the structure and format of this output related to
2180 PLC1 of the iTrust SWaT system is shown in Listing 4.5, where a condition
2181 was specified on the measurement LIT101 and on actuator MV101:

```
2182    aprogram.point:::POINT
2183     P102 == prev_P102
2184     FIT101 >= 0.0
2185     MV101 one of { 0.0, 1.0, 2.0 }
2186     P101 one of { 1.0, 2.0 }
2187     P102 == 1.0
2188     max_LIT101 == 816.0
2189     min_LIT101 == 489.0
2190     slope_LIT101 one of { -1.0, 0.0, 1.0 }
2191     [...]
2192     LIT101 > MV101
```

```

2193     LIT101 > P101
2194     LIT101 > P102
2195     LIT101 < max_LIT101
2196     LIT101 > min_LIT101
2197     [...]
2198     MV101 < min_LIT101
2199     MV101 < trend_LIT101
2200     P101 >= P102
2201     P101 < max_LIT101
2202     [...]
2203 =====
2204    aprogram.point:::POINT;condition="MV101 == 2.0 && LIT101 <
2205     ↪ max_LIT101 - 16 && LIT101 > min_LIT101 + 15"
2206     MV101 == prev_MV101
2207     P102 == slope_LIT101
2208     MV101 == 2.0
2209     FIT101 > MV101
2210     FIT101 > P101
2211     FIT101 > P102
2212     FIT101 > prev_P101
2213     MV101 >= P101
2214     MV101 >= prev_P101
2215     P101 <= prev_P101
2216 =====
2217    aprogram.point:::POINT;condition="not(MV101 == 2.0 &&
2218     ↪ LIT101 < max_LIT101 - 16 && LIT101 > min_LIT101 + 15)
2219     ↪ "
2220     P101 >= prev_P101
2221     Exiting Daikon.

```

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

2222 In the presented framework, the output is simplified to **two sections**:
2223 the *general section* and the section related to the *user-specified condition*. The
2224 section related to the negated condition is eliminated as it is not relevant in
2225 this context and could lead to potential misinterpretation. Moreover, the
2226 relationships between invariants will be emphasized by utilizing **transi-**
2227 **tive closures**: transitive closure of a relation R is another relation, typically
2228 denoted R^+ that adds to R all those elements that, while not necessarily

related directly to each other, can be reached by a *chain* of elements related to each other. In other words, the transitive closure of R is the smallest (in set theory sense) transitive relation R such that $R \subset R^+$ [65].

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:

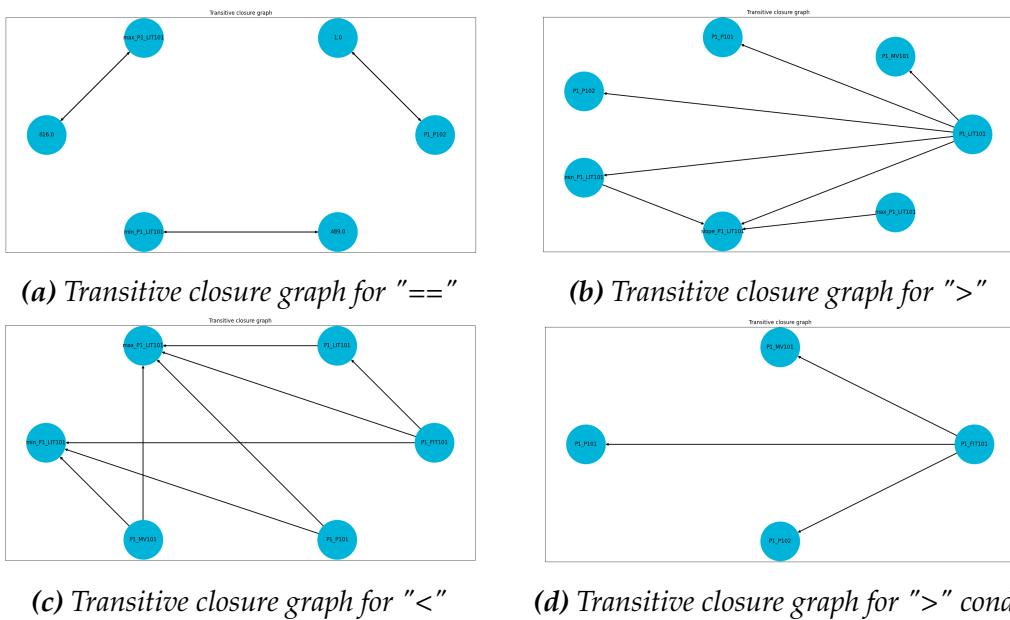


Figure 4.9: Example of transitive closure graphs for invariants in PLC1 of iTrust SWaT system

2244 At the end of this process, still applied to PLC1 of the iTrust SWaT sys-

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

```

2247 1 =====
2248 2 General
2249 3 =====
2250 4 MV101 one of { 0.0, 1.0, 2.0 }
2251 5 P101 one of { 1.0, 2.0 }
2252 6 slope_LIT101 one of { -1.0, 0.0, 1.0 }
2253 7 FIT101 != P101, P102
2254 8 P102 == 1.0
2255 9 max_LIT101 == 816.0
2256 10 min_LIT101 == 489.0
2257 11 LIT101 > MV101
2258 12 LIT101 > P101
2259 13 LIT101 > P102
2260 14 LIT101 > min_LIT101 > slope_LIT101
2261 15 FIT101 >= 0.0
2262 16 P101 >= P102 >= slope_LIT101
2263 17 =====
2264 18 =====
2265 19 MV101 == 2.0 && LIT101 < max_LIT101 - 16 && LIT101 >
2266 20 ↢ min_LIT101 + 15
2267 21 =====
2268 21 slope_LIT101 == P102
2269 22 MV101 == 2.0
2270 23 FIT101 > MV101
2271 24 FIT101 > P101
2272 25 FIT101 > P102
2273 26 MV101 >= P101

```

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

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

2279 **4.2.5.2 Types of Analysis**

2280 In contrast to Ceccato et al.'s solution, which involves manual individual analyses, our proposal is to introduce **two types of semi-automated analysis**.

2283 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.

2287 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.

2293 These two types of analysis will be handled by the Python script `daikonAnalysis.py`, contained in the default directory `$(project_dir)/daikon`.
2294 The script accepts the following command-line arguments:

- 2296 • **-f or --filename**: specifies the enriched dataset, in CSV format, from
2297 which to read the data. The dataset must be located within the directory
2298 containing the enriched datasets specified in the `config.ini` file
2299 (by default `$(project_dir)/daikon/Daikon_Invariants`);
- 2300 • **-s or --simpleanalysis**: performs the analysis on the states of individual
2301 actuators;
- 2302 • **-c or --customanalysis**: performs the analysis on combinations of actual
2303 states of the actuators;
- 2304 • **-u or --uppermargin**: defines a percentage margin on the maximum
2305 value of the measurement;
- 2306 • **-l or --lowermargin**: defines a percentage margin on the minimum
2307 value of the measurement;

2308 The selection of one or both types of analysis is possible. Additionally,
2309 the last two parameters can be used to set a condition on the value of the
2310 measurement. This condition is designed to bypass the transient periods
2311 that occur during actuator state changes and the actual trend changes at
2312 the maximum and minimum values of the measurement.

2313 This approach proves particularly beneficial for the first type of analy-
2314 sis, as it enables more accurate data on the trends of the measurement. By
2315 excluding the transient periods, the analysis can focus on the stable and
2316 meaningful trends, providing improved insights into the behavior of the
2317 system.

2318 **Analysis on single actuator states** Analysis on the states of individual
2319 actuators is the simplest: after the user is prompted to input the measure-
2320 ment, chosen from a list of likely available measurements, the script rec-
2321 ognizes the likely actuators and the relative states of each, using the same
2322 mixed Daikon/Pandas technique adopted in the preliminary analysis dur-
2323 ing the pre-processing phase.

2324 For each actuator and each state it assumes, a single Daikon analysis is
2325 performed, eventually placing the condition on the maximum and mini-
2326 mum level of the measurement.

2327 The result of these analyses are saved in the form of text files in a di-
2328 rectory having the name corresponding to the analyzed actuator and con-
2329 tained in the default parent directory

2330 \$(project_dir)/daikon/Daikon_Invariants/results: each file generated
2331 by the analysis is identified by the name of the actuator, the state and the
2332 condition, if any, on the measurement (see Figure 4.10).

2333 Listing 4.6 provides an illustrative example of the analysis results. In
2334 this case, we focus on the actuator MV101 in state 2.

2335 The condition-generated invariants provide the most interesting in-
2336 sights. For example, at line 21, we observe that when MV101 is set to 2,
2337 the slope of LIT101 is 1, indicating an increasing trend. This leads us to
2338 speculate that MV101 represents the actuator's ON state and is responsible
2339 for filling the tank (LIT101).

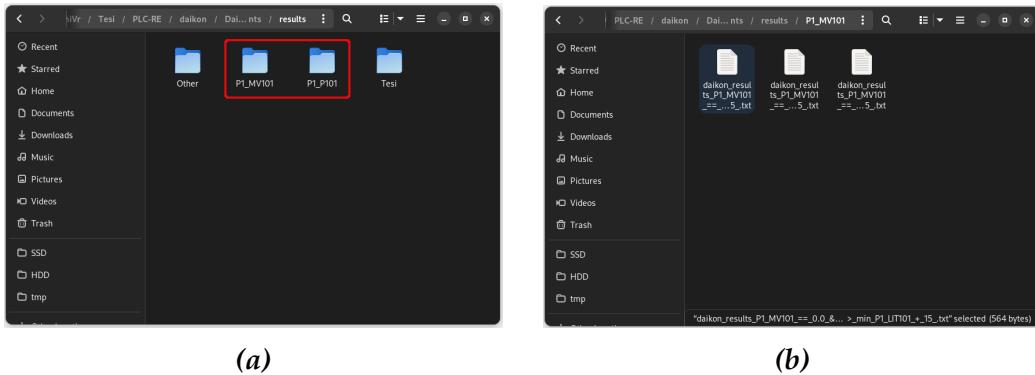


Figure 4.10: Directory (a) and outcome files (b) for the single actuator states analysis

Analysis of the Current System Configuration The analysis of the current system configuration based on the actual states of the actuators is more complex, but at the same time offers more interesting outcomes as it provides better evidence of actuator behavior in relation to the selected measurement.

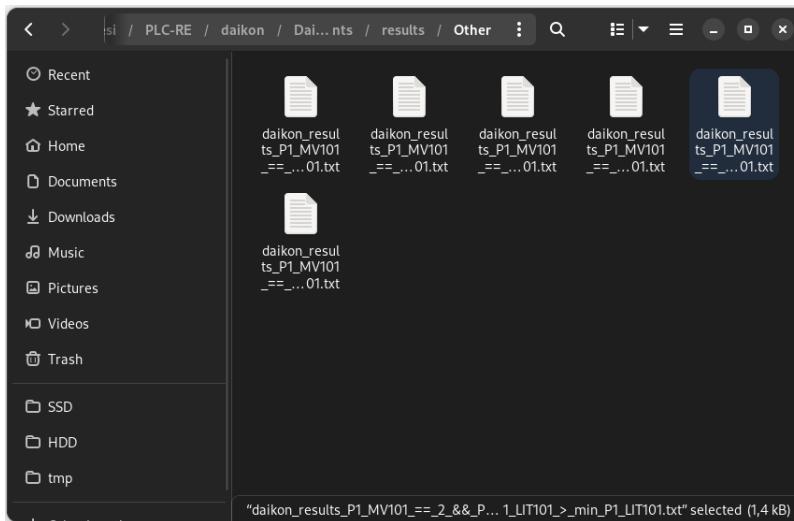


Figure 4.11: Daikon outcome files for system configuration analysis. Each file represents a single system state

2345 For this analysis, the script automatically identifies the configurations of
2346 the actuators that represent different system states (e.g., MV101 == 2, P101
2347 == 1). Daikon analysis is then performed for each of these configurations.

2348 The user is prompted to select the measurement attribute and, if desired,
 2349 specific actuators for studying their configurations. If no actuators are se-
 2350 lected, all previously detected actuators will be considered.

2351 The analysis results are saved in text format within a designated di-
 2352 rectory, located alongside the previous analysis outputs. The filenames of
 2353 these result files follow a specific naming convention based on the analysis
 2354 rule applied (see Figure 4.11).

2355 An example of the obtained outcomes can be seen in Listing 4.7:

```

2356   1 =====
2357   2 General
2358   3 =====
2359   4 MV101 <= P101 ==> FIT101 >= 0.0
2360   5 MV101 <= P101 ==> MV101 one of { 0.0, 1.0, 2.0 }
2361   6 MV101 <= P101 ==> P101 one of { 1.0, 2.0 }
2362   7 MV101 <= P101 ==> slope_LIT101 one of { -1.0, 0.0, 1.0 }
2363   8 MV101 > P101 ==> FIT101 > MV101
2364   9 MV101 > P101 ==> FIT101 > P101
2365  10 MV101 > P101 ==> FIT101 > P102
2366  11 MV101 > P101 ==> FIT101 > slope_LIT101
2367  12 MV101 > P101 ==> MV101 == 2.0
2368  13 MV101 > P101 ==> MV101 > P102
2369  14 MV101 > P101 ==> MV101 > slope_LIT101
2370  15 MV101 > P101 ==> P101 == 1.0
2371  16 MV101 > P101 ==> P101 == P102
2372  17 MV101 > P101 ==> P101 == slope_LIT101
2373  18 MV101 > P101 ==> slope_LIT101 == 1.0
2374  19 [...]
2375  20
2376  21 =====
2377  22 MV101 == 2 && P101 == 1 && LIT101 < max_LIT101 && LIT101 >
2378    ↢ min_LIT101
2379  23 =====
2380  24 slope_LIT101 == P102 == P101 == 1.0
2381  25 MV101 == 2.0
2382  26 FIT101 > MV101

```

2383 27 FIT101 > P101

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

2384 In contrast to Listing 4.6, the current analysis reveals the presence of impli-
2385 cations in the general invariants section, which were previously absent (re-
2386 maining generic invariants omitted for brevity). These implications offer
2387 valuable insights, as demonstrated in this case. For instance, the invariant
2388 stated in line 18 informs us that if the value of MV101 is greater than P101,
2389 then the slope's value is 1, indicating an increasing tank level. This finding
2390 further corroborates our previous understanding that the state MV101 ==
2391 2 represents the ON state for the actuator responsible for filling the tank,
2392 namely LIT101.

2393 **Refining the Analysis** In certain situations, the outcomes provided by
2394 the semi-automated analyses may not meet the user's expectations. For
2395 instance, the clarity of the slope value may be insufficient, or the user may
2396 wish to delve deeper into a specific aspect of the system to uncover ad-
2397 ditional invariants that were not previously identified. In such cases, the
2398 user has the option to conduct a more targeted and specific invariant anal-
2399 ysis using the Python script `runDaikon.py`, which enables precise investi-
2400 gations of the system.

2401 The script, located in the default directory `$(project_dir)/daikon`, can
2402 be executed with three command-line parameters:

- 2403 • **-f or --filename:** specifies the CSV format *enriched* dataset to read
2404 data from. Even in this case, the dataset must be located within the
2405 directory containing the enriched datasets;
- 2406 • **-c or --condition:** specifies the condition for the analysis, which will
2407 be automatically overwritten in the `.spinfo` file. It is possible to spec-
2408 ify more than one condition, but it is strongly recommended to use
2409 the logical operator `&&` to avoid undesired behaviors in Daikon out-
2410 comes;

- 2411 • **-r or --register:** specifies the directory where to save the text file with
2412 the outcomes.

2413 During the execution of this analysis, a single output file is generated,
2414 containing the discovered invariants. The user can specify the directory
2415 where this file should be saved using the `-r` command-line option. By de-
2416 fault, the file is stored in the directory
2417 `$(project_dir)/daikon/Daikon_Invariants/results`. The output file serves
2418 as a record of the identified invariants and can be examined at a later time.

2419 In conclusion, the integration of these two analysis types, along with
2420 the ability to conduct more refined analysis in the future and the enhanced
2421 output format of Daikon, significantly enhances the completeness, clarity,
2422 and effectiveness of this stage compared to the Ceccato et al.'s framework.

2423 4.2.6 Phase 4: Business Process Analysis

2424 We have made significant revisions to the Business Process Analysis
2425 we are presenting. Instead of relying on the previous Java solution and
2426 proprietary process mining software Disco, we have adopted a **new inte-**
2427 **grated solution** created in Python from scratch. The new solution utilizes
2428 the Graphviz libraries to generate the corresponding activity diagram.

2429 In this updated Business Process, greater emphasis is placed on pro-
2430 cess mining related to the physical system. Our goal is to extract as much
2431 information as possible from the dataset, enabling us to promptly visu-
2432 alize the system's behavior and its various states. This approach allows
2433 us to validate the conjectures and hypotheses formulated in the previ-
2434 ous phases, and potentially uncover hidden patterns that were previously
2435 undisclosed.

2436 On the other hand, the aspect related to network communications was
2437 reconsidered to enable operation with multiple protocols, expanding be-
2438 yond the limitations of Modbus.

2439
2440 Now, let's examine the key aspects of this new phase in greater detail.

2441 4.2.6.1 Process Mining of the Physical Process

2442 The mining of the physical process is performed by the Python script
2443 called `processMining.py`, located in the default directory `$(project_dir)/process-mining`.
2444 This script accepts the following parameters from the command line:

- 2445 • **-f or --filename:** specifies the CSV format *timestamped* dataset to be
2446 mined from. The dataset is obtained from the pre-processing stage
2447 and is located in the default directory `$(project_dir)/process-mining/data`;
- 2448 • **-a or --actuators:** specifies one ore more actuators whose combina-
2449 tions of states are to be analyzed. If this parameter is not provided,
2450 all actuators in the subsystem will be considered;
- 2451 • **-s or --sensors:** specifies one or more measurements for which the
2452 trend will be calculated based on actuator state changes. If this pa-
2453 rameter is omitted, all available measurements will be considered;
- 2454 • **-t or --tolerance:** specifies the tolerance to be taken into account dur-
2455 ing the trend calculation;
- 2456 • **-g or --graph:** shows the resulting activity diagram.

2457 The script processes the dataset in a sequential manner, analyzing each
2458 actuator state change. It calculates the duration in seconds of the system
2459 state, as well as the trend and slope of the specified measurement(s). Addi-
2460 tionally, the script stores the next system state and the measurement values
2461 corresponding to the state change points, allowing for the identification of
2462 relative setpoints. These results are gradually collected and stored in a dic-
2463 tionary, where the keys represent the system states based on the actuator
2464 configurations, and the values represent the measured values mentioned
2465 above. The dictionary is then saved as a JSON file, which can be accessed
2466 by the user in the `$(project_dir)/process-mining/data` directory. The
2467 name of the file can be specified in the configuration file `config.ini`.

2468 An example of the JSON file obtained in this step is shown in Listing
2469 4.8: the JSON file showcases the structure of the data obtained during the
2470 process.

```
2471 {
2472     "MV101 == 2, P101 == 2": {
2473         "start_value_LIT101": [
2474             534.9366,
2475             495.4483,
2476             497.9998,
2477             495.9586,
2478             495.8016
2479         ],
2480         "end_value_LIT101": [
2481             536.0356,
2482             532.7384,
2483             541.7273,
2484             534.4656,
2485             540.8245
2486         ],
2487         "slope_LIT101": [
2488             0,
2489             0.015,
2490             0.018,
2491             0.016,
2492             0.018
2493         ],
2494         "trend_LIT101": [
2495             "STBL",
2496             "ASC",
2497             "ASC",
2498             "ASC",
2499             "ASC"
2500         ],
2501         "time": [
2502             119,
2503             2464,
2504             2477,
2505             2452,
2506             2440
```

```

2507     ],
2508     "next_state": [
2509       "MV101 == 2, P101 == 1",
2510       "MV101 == 2, P101 == 1",
2511       "MV101 == 2, P101 == 1",
2512       "MV101 == 2, P101 == 1",
2513       "MV101 == 2, P101 == 1"
2514     ]
2515   },
2516   "MV101 == 2, P101 == 1": {
2517     ...
2518   }
2519   "MV101 == 0, P101 == 1": {
2520     ...
2521   }
2522   ...
2523 }
2524 }
```

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.

2541 Conversely, a high standard deviation indicates that the state may occur
 2542 multiple times within the cycle.

2543 An example of the activity diagram generated by the processMining.py
 2544 script is presented in Figure 4.12. This diagram illustrates the system's
 2545 behavior by depicting transitions between different states. Nodes in the
 2546 diagram represent specific system states, while arrows or edges indicate
 2547 the flow between these states.

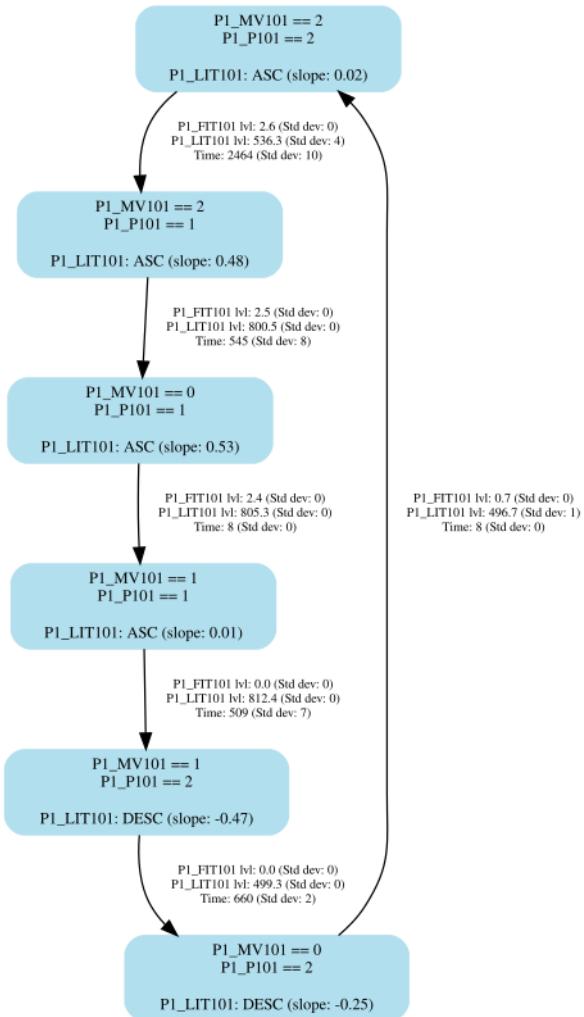


Figure 4.12: Activity diagram for PLC1 of the iTrust SWaT system

2548 The diagram reveals important aspects of the system, such as its **periodic-**

2549 ity and the emphasis on the **temporal sequence of actuator states**. These
2550 insights might not have been clearly evident in earlier stages of the anal-
2551 ysis. The diagram provides a visual representation that allows us to ap-
2552 preciate the recurring patterns and understand how the system evolves
2553 over time. By observing the transitions and relationships between differ-
2554 ent states, we gain a better understanding of the dynamics and behavior
2555 of the system.

2556 It is important to acknowledge that despite considering the tolerance
2557 set for trend and slope calculation, the accuracy of the related data may
2558 vary. For instance, there could be instances where multiple trends exist
2559 within the same node, or a trend may be stable instead of increasing or de-
2560 creasing. These discrepancies arise due to **data perturbations**, which were
2561 discussed in Section 4.2.3.2. Additionally, the noise reduction techniques
2562 seen in the same Section were not employed in this analysis. Therefore, it
2563 becomes the responsibility of the user to correctly interpret the outcomes
2564 of this step, taking into account the information presented in the process
2565 mining diagram and the findings from previous analyses. By consider-
2566 ing these factors together, the user can make informed interpretations and
2567 draw accurate conclusions from the results.

2568 **4.2.6.2 Network Data**

2569 The incorporation of network data into Business Process Analysis has
2570 been reconsidered to shift away from a *single-protocol* solution based on
2571 Modbus. Instead, the focus is on adopting a solution that can handle **mul-**
2572 **tiple protocols**, even at the same time.

2573 The main concept was to develop a new Python script that would ex-
2574 tract and process data from PCAP files obtained through network traffic
2575 sniffing. The intention was to export this data to a CSV file, which would
2576 then be used as input for the process mining script, `processMining.py`,
2577 and integrated with the physical system data to derive commands to the
2578 actuators.

2579 The analysis of the network data revealed that different protocols ex-
 2580 hibit distinct behaviors when commanding changes in actuator states. Con-
 2581 sequently, the previous approach proposed by Ceccato et al. becomes **im-**
 2582 **practical** when attempting to detect system state changes through network
 2583 commands sent to the actuators.

2584 However, considering that system state changes have already been
 2585 identified through the process mining of the physical process, and with
 2586 the corresponding event timestamps available, we propose an alternative
 2587 approach to Ceccato et al.'s method. Instead of seeking the correspon-
 2588 dence between network events and physical process events at the same
 2589 instant, we suggest reversing the perspective. By focusing on the corre-
 2590 spondence between a given event occurring in the physical process at a
 2591 specific moment and the corresponding event in the network data at the
 2592 same moment, it should be possible to achieve a similar, if not superior,
 2593 outcome compared to the previous solution.

2594

4.2.7 Summary

2595 In Table 4.1, we provide a *phase-by-phase* summary of the enhance-
 2596 ments achieved by the proposed framework compared to the limitations
 2597 observed in Ceccato et al., as presented in Table 3.1.

Phase	Enhancements
General Improvements	<ul style="list-style-type: none"> - Independence from the analyzed system achieved 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.

Phase 0: 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.
Phase 1: 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).
Phase 1: 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.

Phase 2: 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.
Phase 3: 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.
Phase 4: 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

2598 HAVING introduced the innovative framework and highlighted its po-
2599 tential in the preceding chapter, we now turn our attention to the
2600 case study where we will apply this framework. As previously mentioned
2601 in Chapter 4 and demonstrated through various examples in the same
2602 chapter, our focus will be on the **iTrust SWaT system** [36], developed by
2603 the iTrust – Center for Research in Cyber Security of University of Singa-
2604 pore for Technology and Design [30]. The acronym SWaT represents *Secure*
2605 *Water Treatment*.

2606 The iTrust SWaT system is a testbed that replicates on a small scale
2607 a real water treatment plant arises to support research in the area of cy-
2608 ber security of industrial control systems and has been operational since
2609 March 2015: it is still being used by students at the University of Singapore
2610 for educational and training purposes and is available to organizations to
2611 train their operators on cyber physical incidents.

2612 5.1 Architecture

2613 In contrast to the virtualized testbed discussed in Section 3.2.1 by Cec-
2614 cato et al., the iTrust SWaT system is composed entirely of physical hard-
2615 ware components. It encompasses various elements, starting from field

2616 devices and extending to PLCs, HMI, SCADA workstations, and the SCADA
2617 server (also referred to as the *historian*). The historian is responsible for
2618 recording data from field devices for further analysis. In the upcoming
2619 sections, we will delve deeper into the architecture of the physical process
2620 and the communication network.

2621 **5.1.1 Physical Process**

2622 The physical process of the SWaT consists of six stages, denoted P1
2623 through P6 [66][67]. These stages are:

- 2624 P1. **taking in raw water:** feeds unfiltered water into the system
- 2625 P2. **chemical dosing:** adds chemicals to water useful for initial pretreat-
2626 ment;
- 2627 P3. **Ultra Filtration (UF) system:** the water is filtered through a semi-
2628 permeable membrane (ultrafiltration membrane) using the liquid pres-
2629 sure, effectively capturing impurities and suspended solids, as well
2630 as removing bacteria, viruses, and other pathogens present in the
2631 water;
- 2632 P4. **dechlorination:** removes residual chlorine from disinfected water
2633 using ultraviolet lamps;
- 2634 P5. **Reverse Osmosis (RO):** performs further filtration of the water;
- 2635 P6. **backwash process:** cleans the membranes in UF using the water pro-
2636 duced by RO.

2637 Figure 5.1 [36] shows a graphical representation of the architecture and the
2638 six stages of the SWaT system.

2639 The SWaT system incorporates an array of sensors that play a crucial
2640 role in monitoring the system's operations and ensuring their safety. These
2641 sensors are responsible for continuously collecting data and providing in-
2642 teresting insights into the functioning of the system [66]. These sensors
2643 are:

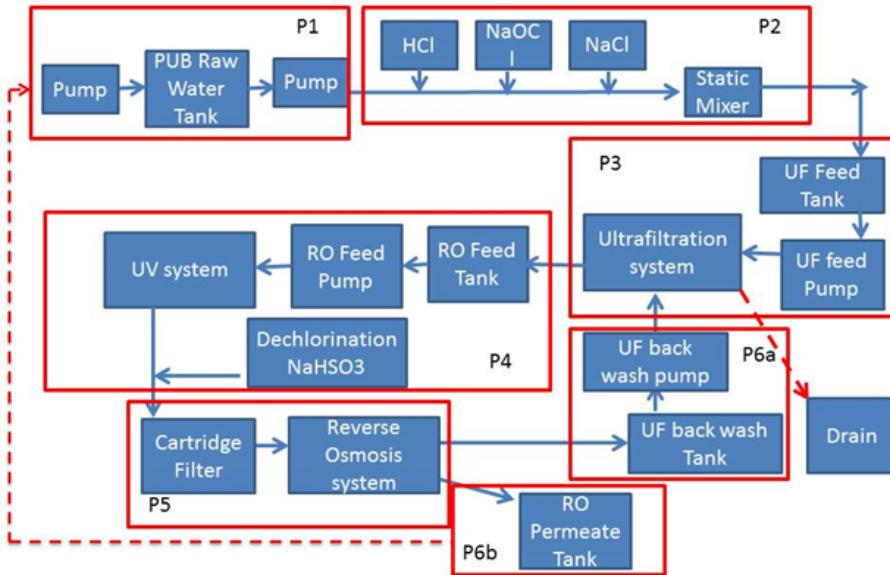


Figure 5.1: SWaT architecture

- Level Indication Transmitter (measured in mm);
- Flow Indication Transmitter (m³/hr);
- Analyser Indicator Transmitter;
 - o Conductivity ($\mu\text{S}/\text{cm}$);
 - o pH;
 - o 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 [66]. 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 [67].

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

2658 5.1.2 Control and Communication Network

2659 The SWaT system's network architecture follows the principles of lay-
 2660 ering and zoning, which enable segmentation and control of traffic within
 2661 the network.

2662

2663 Five layers are present starting from the highest to the lowest:

- 2664 • Layer 3.5 – Demilitarized Zone (DMZ);
- 2665 • Layer 3 – Operation Management (Historian);
- 2666 • Layer 2 – Supervisory Control (Touch Panel, Engineering Worksta-
2667 tion, HMI Control Clients);
- 2668 • Layer 1 – Plant Control Network (PLCs) (Star Network);
- 2669 • Layer 0 – Process (Actuator/Sensors and Input/output modules)
2670 (Ring Network).

2671 PLCs at Layer 1 communicate with their respective sensors and actu-
 2672 ators at Layer 0 through a conventional ring network topology based on

2673 EtherNet/IP, to ensure that the system can tolerate the loss of a single link
2674 without any adverse impact on data or control functionality.

2675 PLCs between the different process stages at Layer 1 communicate
2676 with each other through a star network topology using the CIP protocol
2677 on EtherNet/IP, previously discussed in Section 2.2.6.3.

2678 Regarding zoning, the SWaT system is divided into three zones, each
2679 containing one or more layers. These zones are, in descending order of
2680 security level:

- 2681 • **Plant Control Network, or Control System:** includes layers from 0
2682 to 2;
- 2683 • **DMZ:** includes Layer 3.5;
- 2684 • **Plant Network:** includes Layer 3;

2685 Figure 5.3 [36] provides a clearer visualization of the zoning and layer
2686 division within the network architecture of the SWaT system. This dia-
2687 gram highlights the distinct zones and their corresponding layers, offering
2688 a comprehensive overview of the system's network structure.

2689 A specific IP address is associated with each device: in Table 5.1 we
2690 report the addresses for the PLCs, historian, and Touch Panel in the *Plant*
2691 *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 SWaT

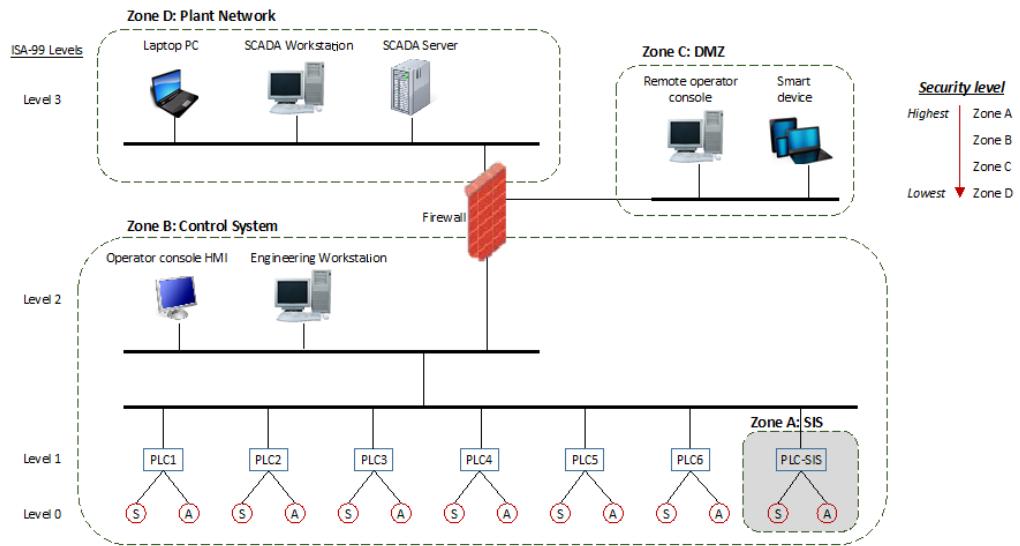


Figure 5.3: SWaT network architecture and zoning

2692 5.2 Datasets

2693 To facilitate the study and testing of technologies related to cyber secu-
 2694 rity in Industrial Control Systems and critical infrastructure, iTrust offers
 2695 researchers worldwide the opportunity to access a range of datasets [68].
 2696 These datasets consist of a collection of data obtained from the SWaT sys-
 2697 tem, encompassing information on both the physical processes and net-
 2698 work communications. The data is organized into different years, and
 2699 researchers can request access to these datasets for their analysis and ex-
 2700 perimentation purposes.

2701 **Physical Process Datasets** The datasets containing information about
 2702 the physical processes are provided in CSV format files. These files en-

2703 compass data collected during different time intervals, which can vary
2704 from a few hours to entire days. The granularity of the data is typically
2705 at a one-second interval, although there may be some exceptions. The col-
2706 lected data primarily consists of timestamped sensor measurements and
2707 actuator status values for each PLC, describing the physical properties of
2708 the testbed in operational mode.

2709 **Network Communications Datasets** Network communications are typi-
2710 cally available in the form of *Packet Capture* format (PCAP) files. These files
2711 contain captures of communication network traffic, allowing researchers
2712 to analyze and examine the network interactions. In some instances, CSV
2713 files are provided instead of PCAP files, featuring different characteristics
2714 for the collected data.

2715 5.2.1 Our Case Study: the 2015 Dataset

2716 The dataset selected as a case study to apply the framework discussed
2717 in the previous chapter is specifically the dataset from the year 2015 [69].
2718 The main reason for this choice is the unique characteristics found in the
2719 physical process dataset that are not present in datasets from subsequent
2720 years.

2721 **Physical Process Data** The data collection process lasted 11 consecutive
2722 days, 24 hours per day. During the first 7 days, the system operated nor-
2723 mally without any recorded attacks. However, attacks were observed dur-
2724 ing the remaining 4 days. The collected data reflects the impact of these
2725 attacks, leading to the creation of two separate CSV files: one containing
2726 the recorded data of SWaT during the system's regular operations, and
2727 the other containing data recorded during the days of the attacks. To en-
2728 sure accurate information about the system, the dataset pertaining to the
2729 normal operations, which spans seven days, was chosen for analysis.

2730 Data collection occurs at a frequency of one data point per second,
2731 with the assumption that significant attacks cannot occur within a shorter

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

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

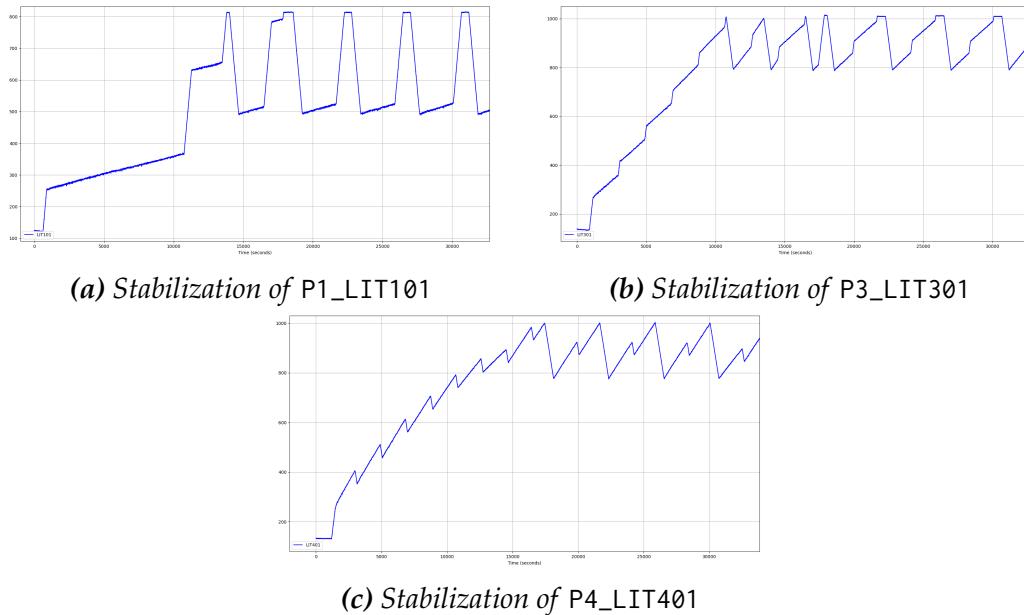


Figure 5.4: SWaT stabilization

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

Network Traffic The network traffic was collected using an appliance from a well-known network hardware manufacturer and was made available only in CSV format and not PCAP format. Table 5.2 shows some of 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

2744 It is important to note that the presence of the string *Modbus* within the
 2745 name of the captured data fields might be misleading, as it may give the
 2746 impression that the communication is taking place via the Modbus proto-
 2747 col rather than the CIP over EtherNet/IP protocol. However, in reality, the
 2748 latter is the actual protocol being used. The indication of Modbus within
 2749 the network traffic dataset could be a result of the appliance used for net-
 2750 work data sniffing, which might have encapsulated CIP over EtherNet/IP
 2751 protocol packets within Modbus frames. This technique is described by
 2752 Snide in [70].

2753 To confirm that the communication protocol used is CIP over Ether-
 2754 Net/IP, we can examine the "Service / Destination Port" field. The destina-
 2755 tion port displayed is TCP port 44818, which corresponds to the TCP port
 2756 commonly used for transporting **explicit messages** in the EtherNet/IP
 2757 protocol (as explained in Section 2.2.6.2). Further evidence supporting the
 2758 use of CIP protocol is observed in the *Modbus Function Code* field, con-
 2759 sistently displaying a value of 76. In the Modbus protocol, there is no

2760 corresponding function associated with that code (as described in Section
2761 2.2.6.1). However, when converting the value 76 from decimal to hex-
2762 adecimal, it translates to **4C**, which corresponds to the **Read Tag Request**
2763 **service in the CIP protocol**. This correlation is further supported by the
2764 presence of the *Application Name* and *Modbus Function Description* fields,
2765 explicitly indicating the protocol name and service within them.

2766 Datasets for years beyond 2015 were not considered due to significant
2767 limitations that impede their usability. For instance, the 2017 dataset ex-
2768 hibits a high level of granularity in the physical system data, combined
2769 with short temporal periods for data gathering. As a result, valuable in-
2770 formation is lost, and the network traffic data cannot be effectively uti-
2771 lized. On the other hand, post-2017 datasets suffer from the issue of being
2772 "spurious." This means that no distinction is made between the data col-
2773 lected during normal operations of the SWaT system and the data associ-
2774 ated with system attacks, as was done in the 2015 dataset. Furthermore,
2775 the 2018 dataset lacks network traffic data altogether.

2776 Regarding network traffic related to the year 2017, the only ones that
2777 seemingly remain unaffected by attacks, the data is divided into large
2778 PCAP files weighing more than 6 GB each. Despite their size, these files
2779 only contain a few minutes of data, which is insufficient to cover even
2780 one complete system cycle. This renders the use of these files impractical,
2781 considering the available resources.

2782 Despite their large quantity, the CSV files associated with network traf-
2783 fic for the year 2015 are comparatively smaller in size, just slightly over 100
2784 MB each. These files cover a more extensive time period, which facilitates
2785 easier management and analysis. Prioritizing the physical process dataset
2786 as crucial, I have selected the year 2015 as a case study, even though the
2787 network data may be incomplete.

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

2788 IN THIS chapter, our main objective is to apply the framework and method-
2789 ology introduced in Chapter 4 to the case study of the iTrust SWaT sys-
2790 tem, as illustrated in Chapter 5. The purpose of this analysis is to assess
2791 the effectiveness and potential of the proposed framework within the con-
2792 text of a system that closely replicates a real-world water treatment plant,
2793 albeit on a smaller scale.

2794 Due to the complexity of the system and the limited space available
2795 in this thesis, we will not conduct a comprehensive analysis and reverse
2796 engineering of the entire system. Instead, we will focus on specific parts
2797 for analysis. We leave it to the reader or those interested in utilizing the
2798 proposed methodology and framework to complete the analysis, should
2799 they choose to do so.

2800 By focusing on selective components and leaving room for further ex-
2801 ploration, we strike a balance between providing valuable insights and
2802 acknowledging the potential for additional research. This approach em-
2803 powers the reader and interested individuals to explore the iTrust SWaT
2804 system further and leverage the proposed methodology and framework
2805 for a more comprehensive analysis.

2806 6.1 Preliminary Operations

2807 Prior to beginning the actual analysis, several preliminary manual op-
2808 erations need to be conducted on the physical process dataset utilized as
2809 a case study, specifically the SWaT system dataset for the year 2015 as out-
2810 lined in Section 5.2.1. To simulate the data-capture process performed by
2811 Ceccato et al. using their scanning tool, the original dataset in XLSX format
2812 (proprietary to Microsoft Excel) was divided into multiple datasets in CSV
2813 format. Each of these datasets corresponds to the individual stages of the
2814 SWaT system and contains the respective registers. These resulting files
2815 were then saved in the directory specified by the `raw_dataset_directory`
2816 directive in the framework configuration file, `config.ini`, ready to be used
2817 in the pre-processing phase. Furthermore, the headers were manually re-
2818 named by adding a prefix from P1_ to P6_ to each register's name. This
2819 prefix indicates the stages, ensuring that each register is easily identifiable
2820 and linked to its corresponding stage.

2821 6.2 Planning the Analysis Strategy

2822 The complexity of the system being analyzed necessitates the adoption
2823 of a deliberate strategy for the analysis. It is not feasible to rely on trial and
2824 error or attempt every possible combination between stages. The former
2825 approach may overlook crucial relationships between PLCs or between
2826 registers, while the latter may result in excessive and unproductive efforts
2827 if the specific portion of the system being analyzed lacks significant infor-
2828 mation or relationships. A sound analysis strategy helps us focus on the
2829 important parts of the system, improving the quality of the analysis and
2830 leading to better process comprehension. By prioritizing our attention, we
2831 can gain a deeper understanding of the crucial components, resulting in
2832 more informed decision-making and a comprehensive understanding of
2833 the overall processes.

2834 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.

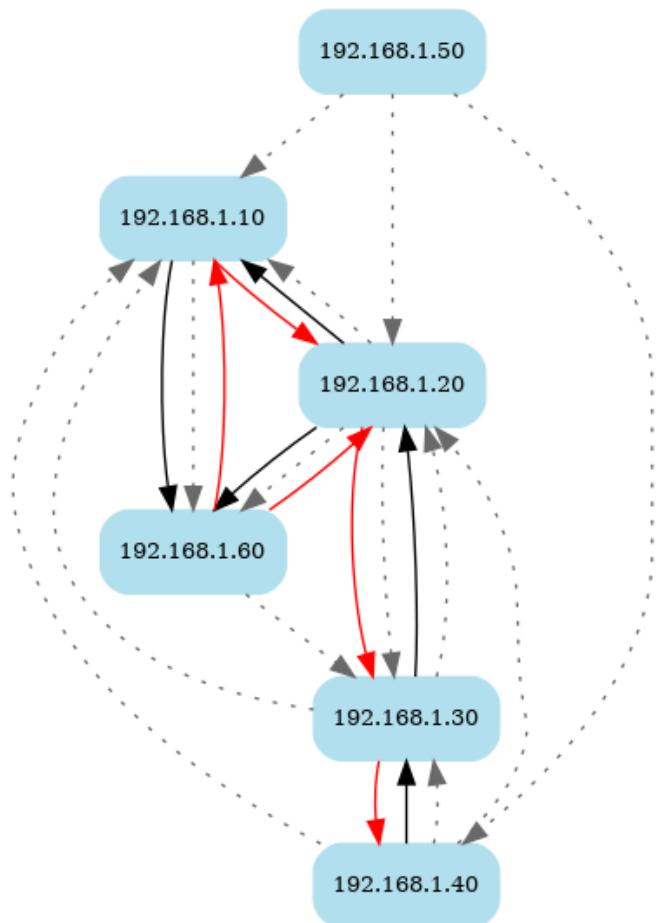


Figure 6.1: Simplified graph of the iTrust SWaT system network

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,

2846 PLC6 communicates with both PLC1 and PLC2. On the other hand, the
2847 gray dotted arrows indicate communications for which we have knowl-
2848 edge of a response, but the corresponding request is unknown. For the
2849 purposes of our analysis strategy, we will not consider these communica-
2850 tions within this context.

2851 Based on our observations, the analysis strategy we will adopt involves
2852 considering sequential pairs of PLCs to effectively capture the relation-
2853 ships and implications between registers.

2854 **6.3 Reverse Engineering of the iTrust SWaT Sys-
2855 tem**

2856 Before we delve into the analysis, it is important to provide some pre-
2857 liminary remarks.

2858 **Analysis structure** Firstly, the analysis will be structured as a schematic
2859 analysis due to space constraints, which prevent us from presenting the
2860 extensive inferences and reasoning regarding the system in full detail.
2861 Therefore, following the analysis strategy outlined in Section 6.2, we will
2862 concentrate exclusively on the pairs of PLCs comprising PLC1-2, PLC2-
2863 3, and PLC3-4. However, the general procedure of the methodology and
2864 how to reason about the data obtained from the framework have already
2865 been demonstrated through examples on the PLC1 of the SWaT system in
2866 Chapter 4. We encourage readers to refer to those examples for a more
2867 comprehensive understanding. In this analysis, our focus will be on illus-
2868 trating the conjectures and properties that arise from the analysis, utilizing
2869 tables and the outputs generated during the analysis.

2870 **Defining subsystem time duration** The second premise addresses the
2871 process of defining the subsystems to be analyzed, which were obtained
2872 during the pre-processing phase, during the merge phase of the individual

2873 datasets. Apart from the projection determined by the considered PLCs,
2874 a time-based selection of the analysis period has also been performed (see
2875 Section 4.2.3.1). This selection spans a duration of 20,000 seconds, which
2876 is equivalent to approximately five and a half hours or roughly five sys-
2877 tem cycles. The analysis begins at 100,000 seconds, which corresponds to
2878 approximately 27 hours from the start of the available data. This deliber-
2879 ate selection aims to exclude the initial transient period during which the
2880 SWaT system is initialized. We believe that this time range is more than
2881 sufficient for accurately defining the characteristics of the SWaT system
2882 components.

2883 **Conventions** The third premise introduces a convention that governs
2884 the naming of PLC registers and will be consistently followed throughout
2885 our analysis. According to this convention, registers with similar names,
2886 such as P1_LIT101 and P3_LIT301, P2_MV201 and P3_MV202, are considered
2887 to belong to the **same category or type of register**. This convention allows
2888 us to establish a relationship or correspondence between registers based
2889 on their naming pattern. By grouping registers with similar names, we can
2890 infer that they serve similar functions or represent similar components in
2891 the system, such as level sensors, tanks, pumps, and so on.

2892 **About the Business Process Analysis** In the end, the Business Process
2893 Analysis will focus solely on the physical process part. This is because the
2894 datasets of network traffic captures provided by iTrust for the year 2015
2895 (as discussed in Section 5.2.1) are **incomplete**. While communications re-
2896 lated to measurements are present, those associated with actuators are en-
2897 tirely missing, as well as additional communications related to other sys-
2898 tem characteristics that we observed in the datasets of subsequent years.
2899 As a result, we were unable to incorporate the network event recognition
2900 component into our Business Process Analysis. To implement this com-
2901 ponent, we would require complete and overlapping network data, along
2902 with a clean physical process dataset not affected by system attacks. Un-
2903 fortunately, none of the available iTrust datasets fulfill these criteria.

2904 6.3.1 Reverse Engineering of PLC1 and PLC2

2905 The initial focus of analysis will be on the pair comprising PLC1 and
 2906 PLC2. Let's delve into the main features of this subsystem by examining
 2907 the outcomes obtained from applying the framework to it.

2908 6.3.1.1 Pre-processing - Preliminary Analysis

2909 **Measurements and Actuators Recognition** Listing 6.1 shows the out-
 2910 comes obtained from automatic recognition of likely measurements and
 2911 actuators:

```
2912 1 Actuators:  

2913 2 P1_MV101 [0.0, 1.0, 2.0]  

2914 3 P1_P101 [1.0, 2.0]  

2915 4 P2_MV201 [0.0, 1.0, 2.0]  

2916 5 P2_P203 [1.0, 2.0]  

2917 6 P2_P205 [1.0, 2.0]  

2918 7  

2919 8 Sensors:  

2920 9 P1_FIT101 {'max_lvl': 2.7, 'min_lvl': 0.0}  

2921 10 P1_LIT101 {'max_lvl': 815.1, 'min_lvl': 489.6}  

2922 11 P2_AIT201 {'max_lvl': 256.5, 'min_lvl': 252.9}  

2923 12 P2_AIT202 {'max_lvl': 8.4, 'min_lvl': 8.3}  

2924 13 P2_AIT203 {'max_lvl': 342.8, 'min_lvl': 320.0}  

2925 14 P2_FIT201 {'max_lvl': 2.5, 'min_lvl': 0.0}  

2926 15  

2927 16 Hardcoded setpoints or spare actuators:  

2928 17 P1_P102 [1.0]  

2929 18 P2_P201 [1.0]  

2930 19 P2_P202 [1.0]  

2931 20 P2_P204 [1.0]  

2932 21 P2_P206 [1.0]
```

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

2933 Based on the results presented in Listing 6.1, the framework has iden-
 2934 tified P1_MV101, P1_P101, P2_MV201, P2_P203, and P2_P205 as **probable ac-**
 2935 **tuators**. The actuators denoted by the *Pxxx* notation are binary actuators

2936 (not boolean!), meaning they have two states represented by the values
2937 1 and 2. Conversely, the actuators identified by the *MVxxx* notation are
2938 ternary actuators with three distinct states: 0, 1, and 2.

2939 To simplify the analysis, we have arbitrarily categorized the registers
2940 identified by the notation *MVxxx* as **valves** and the registers identified by
2941 the notation *Pxxx* as **pumps**. It is important to note that this distinction
2942 is solely for convenience and does *not* necessarily reflect the actual role or
2943 function of these actuators within the system.

2944 P1_FIT101, P1_LIT101, P2_AIT201, P2_AIT202, P2_AIT203, and P2_FIT201
2945 have been identified as **likely measurements**. Upon analyzing the range
2946 of values for register P1_LIT101, we observe a significant difference be-
2947 tween the maximum and minimum values. This observation leads us
2948 to speculate that P1_LIT101 could be identified as a **level sensor for the**
2949 **tank controlled by PLC1**. However, when examining registers P1_FIT101,
2950 P2_FIT201, P2_AIT201, and P2_AIT202, the small difference between their
2951 maximum and minimum values makes it unlikely that they represent ad-
2952 dditional tanks.

2953 Regarding register P2_AIT203, although the range of values is not as
2954 wide as in the case of P1_LIT101, it is still worth examining more closely. It
2955 is possible that P2_AIT203 indicates the presence of a small tank. However,
2956 considering our speculation that the other P2_AIT20x registers are not tank
2957 level sensors, it is uncertain whether P2_AIT203 falls into that category as
2958 well. Further analysis is required to confirm its role within the system.

2959 Some registers have been identified as **hardcoded setpoints or spare**
2960 **actuators** based on their constant values. These registers exhibit similar-
2961 ties to the previously recognized pump registers. It is plausible to spec-
2962 ulate that these registers could correspond to **spare actuators**. Moreover,
2963 the constant value of 1 associated with these registers suggests that it may
2964 represent the **OFF state** of the pumps.

2965 **Actuator State Durations** To gain a deeper understanding of the differ-
2966 ent states (0, 1, and 2) associated with valves P1_MV101 and P2_MV201, we

2967 can analyze the duration of each state. Listing 6.2 provides information
 2968 regarding the duration (in seconds) of states for these specific actuators:

```

2969 1 Actuator state durations:
2970 2 P1_MV101 == 0.0
2971 3 9 9 10 9 9 10 9 9 10 9
2972 4
2973 5 P1_MV101 == 1.0
2974 6 1174 1168 1182 1160 1172
2975 7
2976 8 P1_MV101 == 2.0
2977 9 669 3019 3012 3000 2981
2978 10
2979 11 P2_MV201 == 0.0
2980 12 8 8 8 9 9 8 9 9 9 9
2981 13
2982 14 P2_MV201 == 1.0
2983 15 1057 1057 1045 1038 1039
2984 16
2985 17 P2_MV201 == 2.0
2986 18 120 3135 3144 3127 3109

```

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

2987 It is evident that the duration of **state 0 is relatively short**, averaging
 2988 around 8-10 seconds, while the other states have much longer durations.
 2989 This observation suggests that state 0 of a valve is a **transient state**, in-
 2990 dicating a transitional phase within the valve cycle. However, without
 2991 further information, it is currently not possible to determine the specific
 2992 position of state 0 within the overall valve cycle.

2993 **Actuator State Changes** Now that we have identified P1_LIT101 as the
 2994 supposed level sensor of the tank, we can examine the trend of the tank
 2995 level as the actuators change state. Listing 6.3 provides information on
 2996 the levels of the tank in correlation with the state changes of the P1_P101
 2997 pump:

```

2998 1 Actuator state changes:
2999 2 ...

```

3000	3	P1_LIT101	P1_P101	prev_P1_P101
3001	4	536.0356	1	2
3002	5	533.3272	1	2
3003	6	542.1591	1	2
3004	7	534.8581	1	2
3005	8	540.5890	1	2
3006	9
3007	10	P1_LIT101	P1_P101	prev_P1_P101
3008	11	813.0031	2	1
3009	12	813.0031	2	1
3010	13	811.8256	2	1
3011	14	812.7283	2	1
3012	15	813.3171	2	1
3013	16

Listing 6.3: P1_P101 state changes in relation to P1_LIT101

3014 Based on the speculation that state 1 represents the OFF state of the
 3015 pump and state 2 represents the ON state, we can analyze the data in List-
 3016 ing 6.3. When pump P1_P101 transitions from the ON state to the OFF
 3017 state, the average level of P1_LIT101 is 535. On the other hand, when
 3018 P1_P101 goes from the OFF state to the ON state, the average level of
 3019 P1_LIT101 is 813. These values correspond to the **minimum and maxi-**
 3020 **mum relative setpoints** of P1_P101, respectively.

3021 Based on this information, we can infer that pump P1_P101 is respon-
 3022 sible for **emptying the tank**. Moreover, it can be extended to assume that
 3023 a pump, in general, is responsible for **water outflow**.

3024 Applying the same analysis to the data for valve P1_MV101, which is not
 3025 reported for conciseness, we can speculate that P1_MV101 is responsible for
 3026 **filling the tank**. In this case, states 1 and 2 would represent the **OFF and**
 3027 **ON states of the valve**, respectively. The relative setpoints of P1_MV101
 3028 are approximately 500 (minimum) and 800 (maximum). By extending this
 3029 analysis, we can speculate that a valve, such as P1_MV101, is responsible
 3030 for controlling the **water inflow**.

3031 Regarding the elements controlled by PLC2 and the sensor P1_FIT101,

3032 the analysis does not reveal the presence of another tank. Therefore, we
 3033 cannot determine the exact role of sensors P2_AIT20x, P1_FIT101 and P2_FIT201
 3034 at this point. However, there is a similarity observed between the relative
 3035 setpoints of P1_P101 and those of P2_MV201, P2_P203, and P2_P205. These
 3036 registers exhibit very similar values during state changes, suggesting a
 3037 potential relationship or similar control behavior between them.

3038 **6.3.1.2 Graphs and Statistical Analysis**

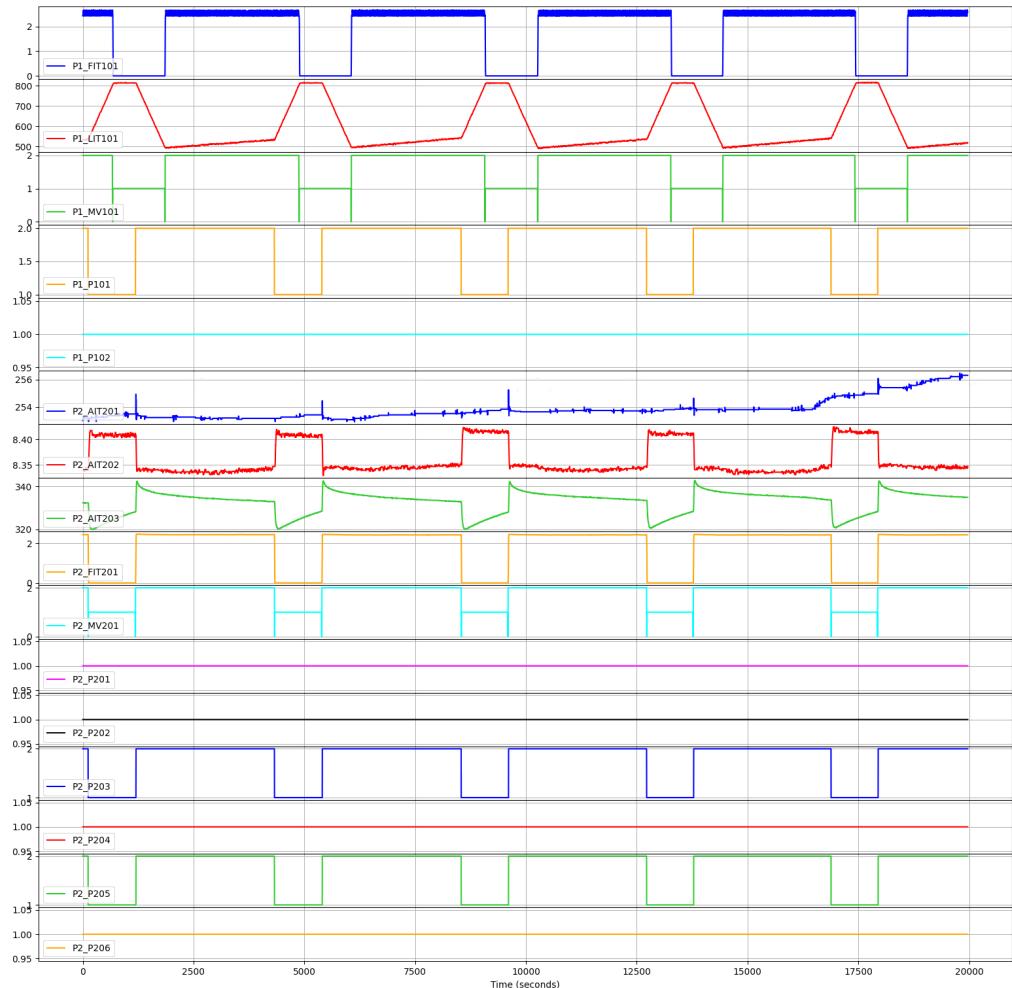


Figure 6.2: Chart of PLC1-2 registers

3039 Figure 6.2 illustrates the graphical representation of the registers in

3040 PLC1 and PLC2 and their respective trends.

3041 The image provides additional support to the conjectures made during
 3042 the preliminary analysis regarding the spare actuators. Furthermore, it is
 3043 evident from the graphs that these spare actuators **do not appear to influ-**
 3044 **ence the trend** of any of the measurements. Therefore, based on this obser-
 3045 vation, we can confidently exclude these registers from further graphical
 3046 analysis.

3047 Figure 6.3 shows a clearer representation of the subsystem after remov-
 3048 ing the spare actuators.

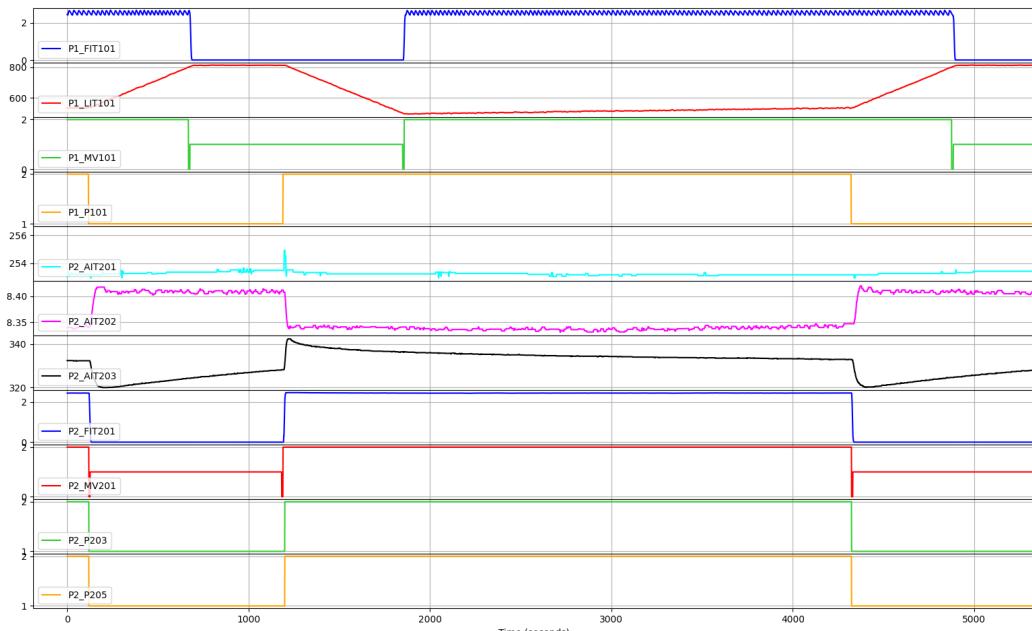


Figure 6.3: Chart of PLC1-2 registers without spare actuators (particular)

3049 Figure provides furthermore additional insights that allow us to spec-
 3050 ulate on aspects that remained unexplained during the preliminary analy-
 3051 sis. The most prominent aspect that stands out is the relationship between
 3052 the level behavior of P1_LIT101 and the states of P1_P101 and P1_MV101. It
 3053 appears evident that these two actuators **do not exhibit complementary**
 3054 **behavior**, meaning that their states do not alternate in an ON-OFF pattern.
 3055 Instead, they can remain in the same state, either ON or OFF, for extended

periods of time. When they are both in the ON state, **the growth of the tank level is slow**. However, as soon as P1_P101 switches to the OFF state, the tank level starts to **increase at a faster rate**. Therefore, when both of these actuators are in the ON state the influx of water into the tank exceeds the outflow of water from it. Conversely, when P1_MV101 switches to the OFF state, the tank level **decreases**. During periods when both actuators are in the OFF state, there is a *plateau* where the tank level remains relatively **stable**.

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

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

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

By examining the trend of valve P2_MV201, an additional speculation 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:

```

3100 1 P2_P206 == P2_P204 == P2_P202 == P2_P201 == P1_P102 == 1.0
3101 2 P2_P205 == P2_P203
3102 3 max_P1_LIT101 == 816.0
3103 4 min_P1_LIT101 == 489.0
3104 5 max_P2_AIT201 == 257.0
3105 6 min_P2_AIT201 == 252.0
3106 7 max_P2_AIT202 == 9.0
3107 8 min_P2_AIT202 == 8.0
3108 9 max_P2_AIT203 == 343.0
3109 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.

3117 **Analysis on Single Actuator States** We proceed with the examination
3118 of the invariants derived from the *first of the two semi-automatic analysis*
3119 discussed in Section 4.2.5.2. Specifically, we will focus on the analysis con-
3120 cerning the states of individual actuators in relation to a specific measure-
3121 ment, which in our case pertains to the tank represented by P1_LIT101. For
3122 illustrative purposes, let's consider states 1 and 2 (OFF e ON respectively)
3123 of valve P1_MV101 as an example (we will disregard state 0 as it is consid-
3124 ered transient). The conditional invariants pertinent to this scenario can
3125 be found in Listing 6.5.

```

3126   1 =====
3127   2 P1_MV101 == 1.0 && P1_LIT101 < max_P1_LIT101 - 16 &&
3128     ↪ P1_LIT101 > min_P1_LIT101 + 15
3129   3 =====
3130   4 ...
3131   5 P2_P205 == P2_P203 == P2_MV201 == P1_P101 == 2.0
3132   6 P1_FIT101 == 0.0
3133   7 slope_P1_LIT101 == -1.0
3134   8 P2_FIT201 > P1_FIT101
3135   9 P2_FIT201 > P1_MV101
3136  10 P2_FIT201 > P1_P101
3137  11 ...
3138  12
3139  13 =====
3140  14 P1_MV101 == 2.0 && P1_LIT101 < max_P1_LIT101 - 16 &&
3141     ↪ P1_LIT101 > min_P1_LIT101 + 15
3142  15 =====
3143  16 slope_P1_LIT101 == P1_P102
3144  17 P1_FIT101 > P1_MV101
3145  18 ...
3146  19 P1_MV101 >= P1_P101
3147  20 P1_MV101 >= P2_MV201
3148  21 P1_MV101 >= P2_P203
3149  22 P2_P203 >= P1_P101
3150  23 ...

```

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

3151 To prevent transient periods caused by water flow stabilization when

3152 the actuators change state, a condition is imposed on the level of P1_LIT101.
3153 However, this condition may result in an incomplete understanding of the
3154 system's behavior. To address this, a manual refinement of the analysis is
3155 required, utilizing the runDaikon.py script as outlined in Section 4.2.5.2.

3156 Based on the analysis, the following observations can be made when
3157 the valve is in the OFF state:

- 3158 • The slope of P1_LIT101, denoted as slope_P1_LIT101, is negative
3159 (line 7). This indicates a **downward trend** in the tank level, as we
3160 have seen in Section 6.3.1.1.
- 3161 • P1_P101 is in state 2, or ON (line 5).
- 3162 • P1_FIT101 is zero (line 6).
- 3163 • P2_FIT201 has a value greater than 2 (line 10).

3164 On the other hand, when the valve is in the ON state:

- 3165 • slope_P1_LIT101 is positive (line 16). This indicates an **upward trend**
3166 in the tank level, as we have seen in Section 6.3.1.1.
- 3167 • P1_FIT101 assumes a value greater than 2. The combination of this
3168 finding, along with the previous one regarding the same register,
3169 strengthens the hypothesis that this is indeed a flow sensor.
- 3170 • P1_P101 can be in either the ON or OFF state, as we have seen in
3171 Section 6.3.2.2.

3172 When conducting a manual analysis using the runDaikon.py script on
3173 tank levels that fall outside the range defined by the previous condition, it
3174 does not yield useful slope information. This situation can occur because,
3175 despite the noise attenuation applied to the tank level sensor data, if there
3176 is even a single cycle in the system where the calculated slope deviates
3177 from the expected outcome, it can adversely affect the entire Daikon anal-
3178 ysis. Figure 6.4 shows this behavior.

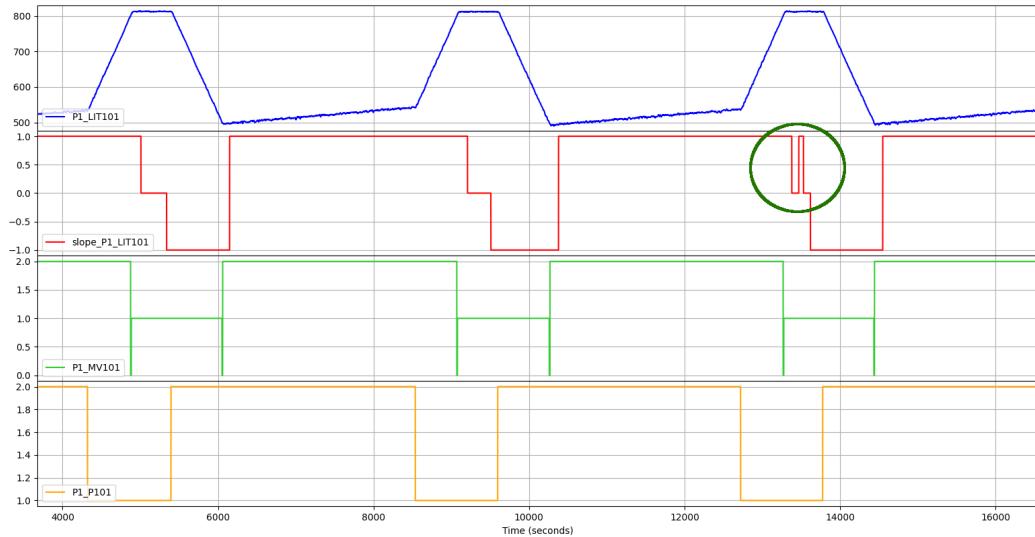


Figure 6.4: Slope calculation anomaly (in the circle)

3179 **Analysis of the Current System Configuration** We conclude our analy-
 3180 sis of invariants by considering the *second semi-automatic analysis* outlined
 3181 in Section 4.2.5.2, which focuses on the effective states of the system. Due
 3182 to the comprehensive nature of this analysis, we will not provide a de-
 3183 tailed report of the outputs to maintain brevity. However, this analy-
 3184 sis confirms the findings observed in the analysis of individual actuator
 3185 states. Additionally, one crucial piece of information becomes apparent
 3186 for future steps: the changing state of the actuators controlled by PLC2 **do**
 3187 **not impact the behavior of the tank** controlled by PLC1. Indeed, it is suf-
 3188 ficient to examine the invariant pertaining to the slope of the tank to verify
 3189 this observation.

3190 6.3.1.4 Business Process Mining and Analysis

3191 As explained in Section 4.2.6.1, the *process mining phase* applied to the
 3192 physical system provides us with an immediate understanding of the sys-
 3193 tem cycle and the chronological sequence of states. It enables us to de-
 3194 termine the duration of time the system remains in a particular state and
 3195 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.

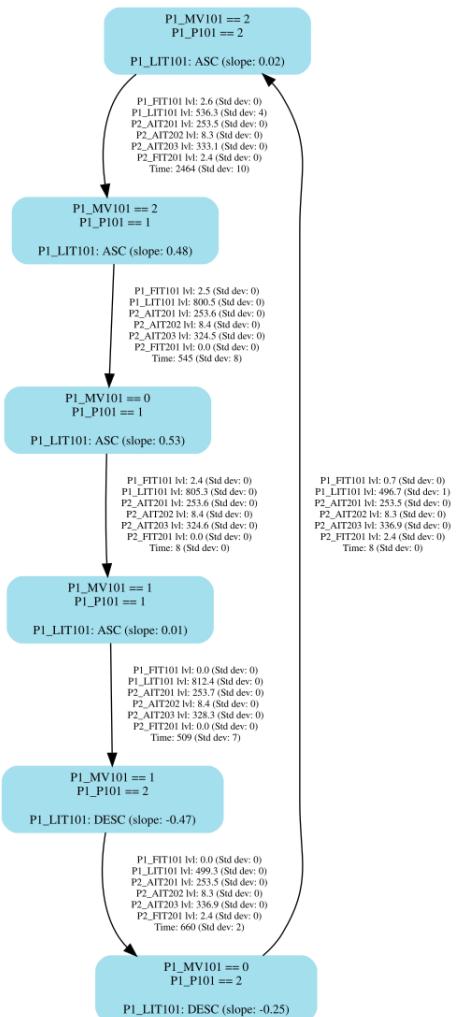


Figure 6.5: Activity diagram for PLC1-2

The activity diagram allows for easy interpretation of the tank level

3204 trend and slope within different states. The states where the valve has a
3205 value of 0 can be disregarded due to their short duration. Similar to the
3206 graphical analysis, there is an observed change in slope during the increasing
3207 trend of the tank level between the states [P1_MV101 == 2, P1_P101
3208 == 2] and [P1_MV101 == 2, P1_P101 == 1], where the slope changes from
3209 0.02 to 0.48.

3210 Additionally, the timing analysis on the edges reveals that the tank
3211 takes longer to fill than to empty, and the system remains in the [P1_MV101
3212 == 1, P1_P101 == 1] state for approximately 8 minutes (509 seconds).

3213 Regarding the state [P1_MV101 == 1, P1_P101 == 1], there appears
3214 to be a discrepancy between the trend of the tank level as reported in the
3215 activity diagram and the conjectures made in the previous phases of the
3216 analysis. The activity diagram correctly depicts an increasing trend in the
3217 tank level between the end of the state [P1_MV101 == 2, P1_P101 == 1]
3218 and the end of the state [P1_MV101 == 1, P1_P101 == 1]. However, the
3219 discrepancy arises due to the fact that the interruption of flow during the
3220 valve state change from state ON to state OFF is not immediate, mainly
3221 due to the presence of the transient state 0. Additionally, water continues
3222 to flow within the piping towards the tank for a short duration even after
3223 the valve is closed. After this period, which usually lasts a few seconds,
3224 the tank level stabilizes, and there is no further inflow or outflow of water.

3225 By adjusting the tolerance parameter -t in the processMining.py script,
3226 it is possible to obtain accurate data regarding the behavior of the state
3227 corresponding to the *plateau* observed in the graphical analysis.

3228 The data presented on the arcs in the activity diagram represents the
3229 measurement values relative to the time of system state changes, specifically
3230 the relative setpoints. These values are calculated based on the average
3231 data collected in each cycle. By analyzing this data, we can observe
3232 the tank level values at which the system undergoes configuration changes
3233 and how the trend of the tank level changes accordingly. Specifically, we
3234 can see that the trend transitions from ascending to stable at a tank level
3235 value of 800, from stable to descending at approximately 812, and from

3236 descending back to ascending at around 499. The change in the speed of
3237 tank filling occurs at approximately 535.

3238 Furthermore, the data provides additional support for the hypothesis
3239 that the measurements associated with registers P2_AIT20x are not influ-
3240 enced by the tank's trends.

3241 **6.3.1.5 Properties**

3242 From the conjectures derived from the four phases of the analysis we
3243 will derive the properties of the subsystem we are studying, which will be
3244 placed within a **summary table**. This table contains an integer identifying
3245 the property, the statement of the property itself, and from which of the
3246 four phases it was derived.

#	Statement	Derived from
P1	The registers P1_LIT101, P1_FIT101, P2_AIT201, P2_AIT202, P2_AIT203, and P2_FIT201 are considered measurements.	Preliminary Analysis Graphical Analysis
P2	The registers P1_MV101, P1_P101, P2_MV201, P2_P203, and P2_P205 are considered actuators.	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 represents 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 considered 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 represents 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 transitions 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 transitions 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 Graphical Analysis
P2_AIT203 are directly associated with the op- Invariant Analysis
eration of the pumps.

Table 6.1: Properties of the PLC1-2 subsystem

3247 6.3.2 Reverse Engineering of PLC2 and PLC3

3248 Continuing our analysis of the iTrust SWaT system, our current focus
3249 will be on the registers of PLC3 and any potential relationships they may
3250 have with the registers of PLC2.

3251 6.3.2.1 Pre-processing - Preliminary Analysis

3252 **Measurements and Actuators Recognition** Listing 6.6 shows the out-
3253 comes obtained from automatic recognition of likely measurements and
3254 actuators. After previously identifying the measurements and actuators
3255 of PLC2 in Section 6.3.1.1, the listing exclusively showcases the registers
3256 associated with PLC3.

```

3257 1 Actuators:
3258 2 ...
3259 3 P3_MV301 [0.0, 1.0, 2.0]
3260 4 P3_MV302 [0.0, 1.0, 2.0]
3261 5 P3_MV303 [0.0, 1.0, 2.0]
3262 6 P3_MV304 [0.0, 1.0, 2.0]
3263 7 P3_P302 [1.0, 2.0]
3264 8
3265 9 Sensors:
3266 10 ...
3267 11 P3_DPIT301 {'max_lvl': 20.4, 'min_lvl': 0.0}
3268 12 P3_FIT301 {'max_lvl': 2.4, 'min_lvl': 0.0}
3269 13 P3_LIT301 {'max_lvl': 1014.5, 'min_lvl': 786.5}
3270 14
3271 15 Hardcoded setpoints or spare actuators:
3272 16 ...

```

3273 17 P3_P301 [1.0]

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

3274 From the provided listing, it is evident that the **likely measurements**
 3275 related to PLC3 are P3_DPIT301, P3_FIT301, and P3_LIT301. Drawing an
 3276 analogy with the derived properties from Table 6.1, P3_LIT301 can be as-
 3277 sociated with a **tank level sensor**, while P3_FIT301 may be linked to a
 3278 flow or pressure sensor. However, the specific role of P3_DPIT301 cannot
 3279 be speculated upon at this time.

3280 Regarding the **likely actuators**, they include P3_MV301, P3_MV302, P3_MV303,
 3281 P3_MV304, and P3_P302. By drawing parallels with the previous analysis
 3282 outcomes, it can be inferred that registers P3_MV30x represent valves, while
 3283 P3_P302 corresponds to a pump.

3284 Lastly, there is a **spare actuator** identified as P3_P301. Similar to the
 3285 previous analysis, this is an inactive pump, indicated by the constant value
 3286 of 1 in this register, signifying that it is in the OFF state.

3287 **Actuator State Durations** Let's proceed with the analysis of the actua-
 3288 tor states' duration, as displayed in Listing 6.7. In this analysis, our focus
 3289 will not be on examining the correspondence between values and the ac-
 3290 tual actuator states, as in Section 6.3.1.1. Instead, we will explore whether
 3291 these actuators exhibit any distinct patterns or behaviors based on their
 3292 duration.

```
3293 1 Actuator state durations:  

  3294 2 ...  

  3295 3 P3_MV301 == 1.0  

  3296 4 2527 4154 4154 4154 4094  

  3297 5  

  3298 6 P3_MV301 == 2.0  

  3299 7 36 35 36 35 34  

  3300 8  

  3301 9 P3_MV302 == 1.0  

  3302 10 662 138 654 138 656 139 658 137 656 137
```

```

3303 11
3304 12 P3_MV302 == 2.0
3305 13 62 1783 1596 1787 1591 1791 1576 1803 1540 1782
3306 14
3307 15 P3_MV303 == 1.0
3308 16 2526 4089 4088 4089 4028
3309 17
3310 18 P3_MV303 == 2.0
3311 19 97 96 97 96 96
3312 20
3313 21 P3_MV304 == 1.0
3314 22 689 1832 2206 1838 2203 1840 2191 1852 2152 1831
3315 23
3316 24 P3_MV304 == 2.0
3317 25 43 87 42 89 43 88 43 88 43 88
3318 26
3319 27 P3_P302 == 1.0
3320 28 637 115 632 115 632 114 634 114 632 115
3321 29
3322 30 P3_P302 == 2.0
3323 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 similar be-

3337 havior to P3_MV302, with a cycle number of 10. Furthermore, the durations
 3338 of the ON and OFF states for both actuators appear to be overlapping. This
 3339 suggests a **potential relationship between the two actuators**. Addition-
 3340 ally, considering that P3_P302 is the only pump in PLC3, it is reasonable to
 3341 speculate that it may have an influence on the tank level. Further inves-
 3342 tigation is necessary to validate this speculation and explore the precise
 3343 nature of the relationship between P3_P302 and P3_MV302.

3344 **Actuator State Changes** Based on the analysis of the probable measure-
 3345 ments, we have identified the likely tank, represented by register P3_LIT301.
 3346 In the previous analysis of the PLC1-2 subsystem in Section 6.3.1, the role
 3347 of valve P2_MV201 remained unresolved. We speculated that this actuator
 3348 might be responsible for the incoming water flow to an element outside
 3349 the analyzed subsystem. To investigate this further, we can examine the
 3350 relationship between P2_MV201 and the tank within this subsystem. By
 3351 extracting information from the corresponding setpoints, we can gather
 3352 insights to verify if our speculation holds true.
 3353 Listing 6.8 displays the setpoints associated with the state change of P2_MV201
 3354 in relation to the level of tank P3_LIT301.

	Actuator state changes:		
	...		
	P2_MV201	prev_P2_MV201	P3_LIT301
3355 1	0	2	1000.2240
3356 2	0	1	799.1140
3357 3	0	2	1001.5060
3358 4	0	1	799.1942
3359 5	0	2	1001.5460
3360 6	0	1	799.1140
3361 7	0
3362 8	0
3363 9	0
3364 10

Listing 6.8: P2_MV201 state changes in relation to P3_LIT301

3365 The setpoints provided in the listing support our conjecture. The max-
 3366 imum relative setpoint of 1000 and the minimum relative setpoint very
 3367 close to 800 indicate a correlation that appears intentional. Based on this

3368 information, we can speculate that P2_MV201 is indeed the **valve responsible**
3369 **for the tank** associated with P3_LIT301.

3370 Regarding further information obtained from this step, there is limited
3371 insight available. While P3_P302 appears to be the pump responsible for
3372 emptying the tank associated with P3_LIT301, the analysis of the actua-
3373 tor lifetime indicates twice as many values compared to other actuators.
3374 The pump exhibits setpoint values of approximately 850 and 970 for the
3375 transition from ON to OFF, and 900 and 1000 for the transition from OFF
3376 to ON. On the other hand, P3_MV302 shows numerous state changes and
3377 shares setpoints with values close to 850, 970, 1000, and 900 for the same
3378 transitions. These values align perfectly with those of the P3_P302 pump,
3379 suggesting a potential relationship between the two actuators.

3380 Obtaining information about the remaining registers is challenging at
3381 this stage. Further analysis steps are required to gather additional insights
3382 and uncover more information about the system.

3383 6.3.2.2 Graphs and Statistical Analysis

3384 Graphical analysis can provide valuable insights and help validate the
3385 conjectures made during the preliminary analysis, as well as uncover con-
3386 nections between registers that were not identified in the previous step.
3387 To begin, we will test the hypothesis that valve P2_MV201 is responsible for
3388 filling the tank associated with P3_LIT101. Figure 6.6 displays these regis-
3389 ters, along with other registers whose behavior and relationships with the
3390 tank level we will explore in an attempt to gain a deeper understanding.

3391 The figure shows the particular behavior of P3_LIT103, with two slope
3392 changes during the tank filling period. These slope changes correspond
3393 approximately to the relative setpoints found for P3_MV302 and P3_P302.
3394 However, we will analyze this aspect later. What we can see in relation
3395 to the initial conjecture about the role of P2_MV201 is that the period dur-
3396 ing which the valve remains in the ON state corresponds exactly to the
3397 increasing trend of P3_LIT301. The conjecture thus finds further support.

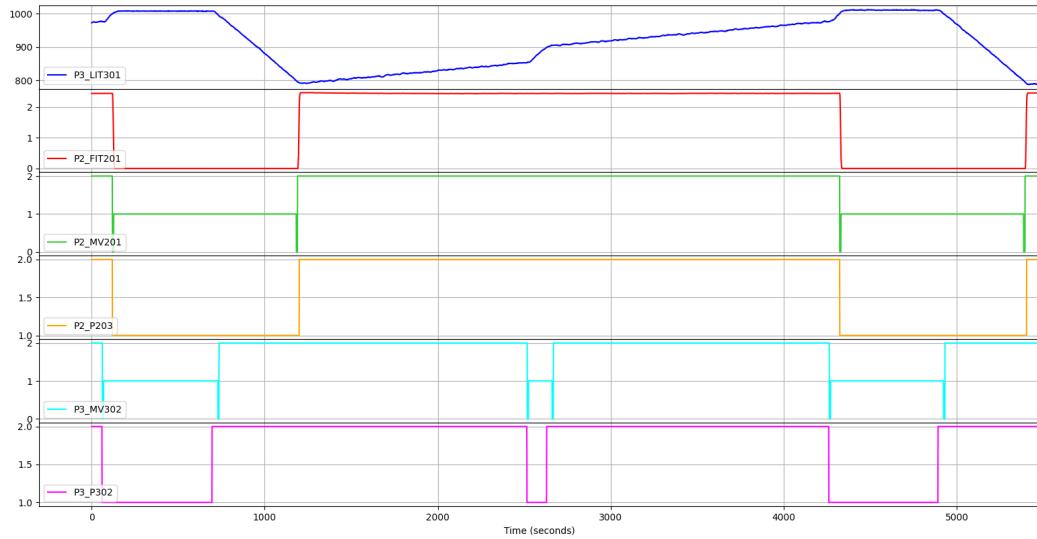


Figure 6.6: Verifying the conjecture about valve P2_MV201

3398 We also note how P2_FIT201 is related to the trend of P2_MV201 and the
 3399 increasing trend of P3_LIT301.

3400 Now let's examine the relationships between the tank represented by
 3401 P3_LIT301 and the other actuators of PLC3 through Figure 6.7.

3402 From the charts, it is evident that the trends of P3_DPIT301 and P3_FIT301
 3403 exhibit similarities. Furthermore, their overall pattern closely follows that
 3404 of the valve P3_MV302. Based on these observations, we can speculate that
 3405 there is a relationship between these registers or that they serve similar
 3406 functions, possibly as **pressure or flow sensors**.

3407 Regarding pump P3_P302, we observe that its OFF state coincides with
 3408 the increasing slope of P3_LIT301 during its upward trend and the entire
 3409 phase when the level remains relatively stable. Conversely, its ON state
 3410 corresponds to the gradual increase and decrease of the water level in the
 3411 tank. This observation provides further evidence to support the hypothe-
 3412 sis that P3_P302 is responsible for emptying the tank.

3413 Valve P3_MV302 exhibits a similar pattern to pump P3_P302. Building
 3414 upon our previous findings, we can speculate that, similar to P2_MV201 in
 3415 the previous analyzed subsystem, P3_MV302 is responsible for controlling

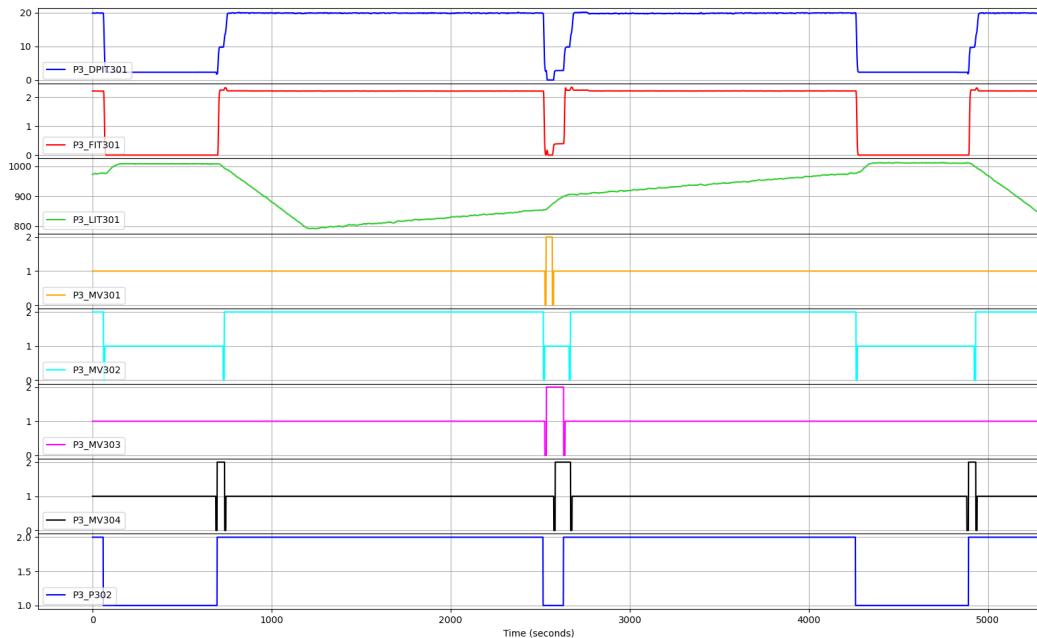


Figure 6.7: PLC3 registers

3416 the incoming flow to another element outside the analyzed subsystem.

3417 Let us now analyze the potential roles of valves P3_MV301, P3_MV303,
 3418 and P3_MV304 in relation to the indicated tank level. It seems unlikely
 3419 that P3_MV304 has any direct impact on the tank level. The valve is acti-
 3420 vated twice within one system cycle, but there are no noticeable changes
 3421 in P3_LIT301 during its first opening. This suggests that the second open-
 3422 ing is also insignificant in terms of tank level. However, it is worth noting
 3423 the slight peaks in P3_FIT301 that occur shortly after the valve's opening.

3424 Regarding the remaining valves, P3_MV301 and P3_MV303, their impact
 3425 on the tank level is still unclear, particularly during the increased slope of
 3426 P3_LIT101 around the second 2500. Figure 6.8 provides us with a compre-
 3427 hensive overview of the situation.

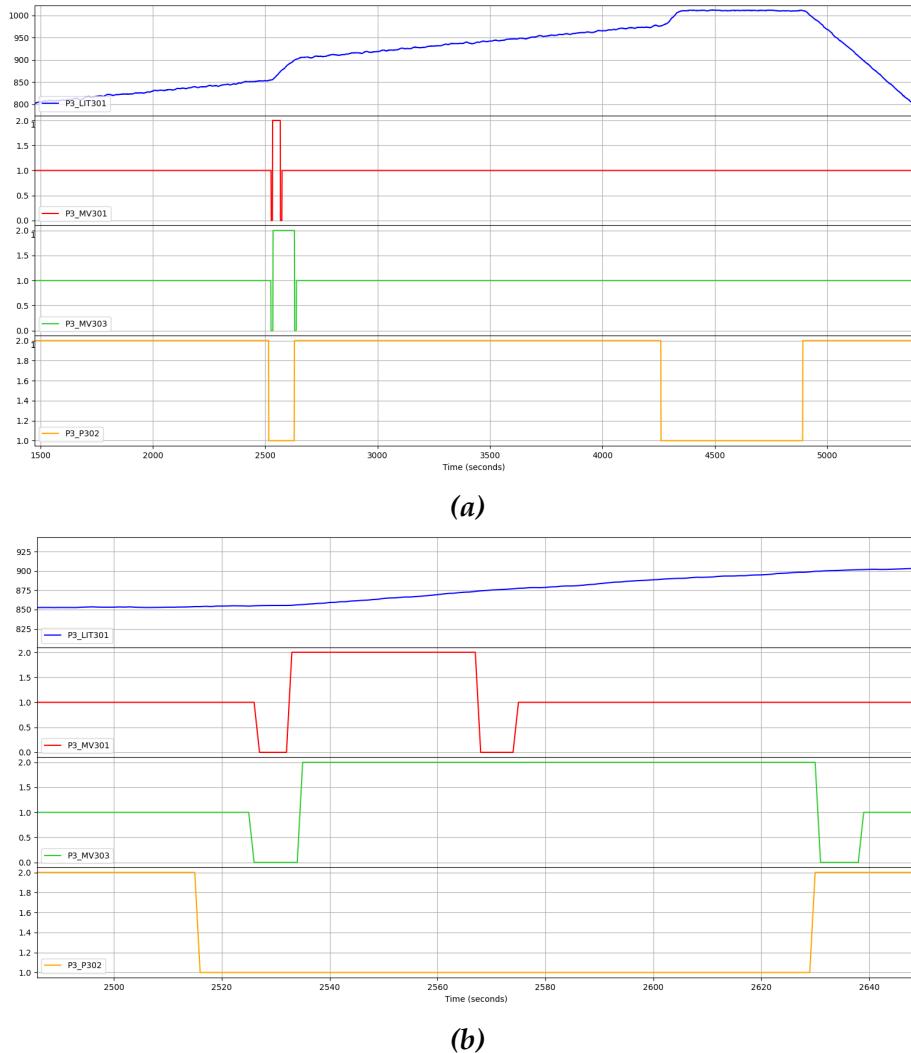


Figure 6.8: P3_MV301 and P3_MV303 analysis

3428 The analysis of the valves P3_MV301 and P3_MV303 indeed presents some
 3429 challenges. However, upon closer examination of Figure 6.8b, we observe
 3430 that when these valves are in the ON state, the slope of P3_LIT301 remains
 3431 relatively constant. This is unexpected, as one would anticipate a steeper
 3432 slope due to the inflow of liquid. Additionally, in Figure 6.8a, we can ob-
 3433 serve that when these valves are in the OFF state, the slope of P3_LIT301
 3434 from the second 4300 remains similar to the section we are currently ana-
 3435 lyzing. Based on these observations, we speculate that these valves either

3436 **do not affect the water level of the tank or have a minimal, undetectable
3437 impact.**

3438 Indeed, Figure 6.7 provides additional insights into the relationship
3439 between the valves P3_MV301, P3_MV303, and P3_MV304 and the sensors
3440 P3_DPIT301 and P3_FIT301. The graph shows that the values of P3_DPIT301
3441 and P3_FIT301 undergo noticeable changes during the activation of these
3442 valves, suggesting a potential connection between them. This observation
3443 supports the hypothesis that the valves and sensors are linked in some
3444 way,

3445 **6.3.2.3 Invariant Inference and Analysis**

3446 **General Invariants** The analysis of general invariants does not yield any
3447 significant insights. However, it confirms the maximum and minimum
3448 values of the measurements seen in Section 6.3.2.1 and identifies the pres-
3449 ence of the spare actuators P3_P301.

3450 **Analysis on Single Actuator States** The analysis of single actuator states
3451 reveals additional information. The resulting invariants provide further
3452 support for the conjecture regarding the roles of P2_MV201 and P3_P302 in
3453 regulating the water level in the tank, with the former responsible for fill-
3454 ing and the latter for emptying. Listing 6.9 presents the specific invariants
3455 involved in this analysis.

```
3456 1 =====
3457 2 P2_MV201 == 1.0 && P3_LIT301 < max_P3_LIT301 - 20 &&
3458 3   ↪ P3_LIT301 > min_P3_LIT301 + 24
3459 4 =====
3460 5 P3_P301 == P3_MV303 == P3_MV301 == P2_P205 == P2_P203 ==
3461 6   ↪ P2_MV201 == 1.0
3462 7 P3_P302 == 2.0
3463 8 slope_P3_LIT301 == -1.0
3464 9 P3_FIT301 > P2_MV201
3465 10 P3_DPIT301 > P3_FIT301 > P3_P302
3466 11 ...
```

```

3467   10 =====
3468   11 P2_MV201 == 2.0 && P3_LIT301 < max_P3_LIT301 - 20 &&
3469      ↪ P3_LIT301 > min_P3_LIT301 + 24
3470   12 =====
3471   13 P2_P205 == P2_P203 == P2_MV201 == 2.0
3472   14 slope_P3_LIT301 == P2_P201
3473   15 P2_FIT201 > P2_MV201
3474   16 P2_FIT201 > P2_P201
3475   17 P2_FIT201 > P3_FIT301
3476   18 ...
3477   19 =====
3478   20 P3_P302 == 1.0 && P3_LIT301 < max_P3_LIT301 - 20 &&
3479      ↪ P3_LIT301 > min_P3_LIT301 + 24
3480   21 =====
3481   22 P2_P205 == P2_P203 == P2_MV201 == 2.0
3482   23 slope_P3_LIT301 == P3_P302 == P2_P201
3483   24 P2_FIT201 > P2_MV201
3484   25 P2_FIT201 > P3_FIT301
3485   26 ...
3486   27 =====
3487   28 P3_P302 == 2.0 && P3_LIT301 < max_P3_LIT301 - 20 &&
3488      ↪ P3_LIT301 > min_P3_LIT301 + 24
3489   29 =====
3490   30 P2_MV201 one of { 1.0, 2.0 }
3491   31 slope_P3_LIT301 one of { -1.0, 1.0 }
3492   32 P3_DPIT301 > P3_P302 > slope_P3_LIT301
3493   33 ...

```

Listing 6.9: Conditional Invariants for P2_MV201 and P3_P302

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

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

3502 However, one observation can be made regarding P3_MV304. Daikon's
3503 analysis reveals two different slopes (increasing and decreasing) when this
3504 valve is in the ON state, which aligns with the observations made in the
3505 Graphical Analysis in Section 6.3.2.2. This finding strengthens the conjecture
3506 that P3_MV304 does not play a significant role in the tank cycle represented
3507 by P3_LIT301.

3508 **Analysis of the Current System Configuration** To simplify the analysis
3509 and facilitate the interpretation of outcomes, two separate groups of actuators
3510 were analyzed. The first group consists of P2_MV201 and P3_P302,
3511 which are conjectured to regulate the level of the tank. The second group
3512 includes P3_MV301, P3_MV303, and P3_MV304, for which it is speculated that
3513 they do not play a role in regulating the tank level. This grouping allows
3514 for a clearer examination of the states of the system and provides an opportunity
3515 to gain a better understanding of the behavior of these actuators.
3516 By analyzing these two groups separately, it becomes easier to draw conclusions
3517 and make comparisons between the different sets of actuators.

3518 The analysis of the first group of actuators, specifically P2_MV201 and
3519 P3_P302, further supports the conjectures made regarding their behavior
3520 in relation to the trend of the tank level. It was necessary to refine the
3521 analysis manually using the runDaikon.py script to obtain more detailed
3522 insights, particularly for the state [P2_MV201 == 2, P3_P302 == 2].

3523 Regarding the second group of actuators, the analysis of their states
3524 only reinforces the hypothesis that their activation does not have a significant
3525 impact on the trend of tank level represented by P3_LIT301.

3526 Unfortunately, no further useful information can be derived from this
3527 phase of the analysis. To address the remaining questions and clarify any
3528 outstanding issues, we will proceed to the next and final phase of the sub-
3529 system analysis.

3530 **6.3.2.4 Business Process Mining and Analysis**

3531 One of the hypotheses that needed to be tested was whether the valves
3532 P3_MV301, P3_MV303, and P3_MV304 have an impact on the tank level in the
3533 section between setpoints 850 and 900. Our initial assumption was that
3534 these actuators do not affect the level detected by P3_LIT301 because, as
3535 observed in the Graphical Analysis, the slope in that interval is similar
3536 to the slope between setpoints 970 and 1000, where these valves are not
3537 involved.

3538 To verify this conjecture, we utilized the JSON file generated by the
3539 processMining.py script. This script allows us to isolate the specific inter-
3540 vals and calculate the slope for each of them. In Listing 6.10, we present
3541 the outcomes of these calculations, providing further insights into the be-
3542 havior of the valves during those intervals.

```
3543 1 "slope_P3_LIT301": [  

3544 2 0.383, # from 970 to 1000  

3545 3 0.395, # from 850 to 900  

3546 4 0.384,  

3547 5 0.395,  

3548 6 0.354,  

3549 7 0.38,  

3550 8 0.388,  

3551 9 0.381,  

3552 10 0.385,  

3553 11 0.386  

3554 12 ] ,
```

*3555 Listing 6.10: Slope calculation of P3_LIT301 for the 850-900 and 970-1000 intervals
3556 related to tank levels*

3555 The provided data in Listing 6.10 illustrates the calculated slopes for
3556 the intervals between 970 and 1000 (odd-numbered lines) and the inter-
3557 vals between 850 and 900 (even-numbered lines). Upon examination, we
3558 observe that the values for these two intervals are nearly identical, ac-
3559 counting for some expected fluctuations. This finding reinforces our initial
3560 conjecture that the P3_MV301, P3_MV303 and P3_MV304 valves do not play a

3561 role in the process of filling and emptying the tank. It is indeed possible to
3562 speculate that the P3_MV30x valves, similar to P3_MV302, might have a role
3563 in a different part of the system that has not been analyzed in the current
3564 context.

3565 Based on the activity diagram obtained from the analysis of actuators
3566 P2_MV201 and P3_P302, the behavior of subsystem PLC2-3 can be summa-
3567 rized as follows:

- 3568 • When the system is in the states [P2_MV201 == 2, P3_P302 == 2] and
3569 [P2_MV201 == 2, P3_P302 == 1], the level of the tank represented by
3570 P3_LIT301 is increasing. The tank level exhibits a faster growth rate
3571 in the latter state.
- 3572 • When the system is in the state [P2_MV201 == 1, P3_P302 == 1], the
3573 tank level remains stable.
- 3574 • When the system is in the state [P2_MV201 == 1, P3_P302 == 2], the
3575 tank level is decreasing.

3576 The **absolute setpoints** for the tank level in subsystem PLC2-3 are de-
3577 fined as follows: the minimum setpoint is 800, the maximum setpoint is
3578 1000, and there are additional setpoints at 850, 900, and 970, which corre-
3579 spond to specific changes in the slope of the tank level.

3580 Taking a closer look at the behavior of sensors P2_FIT201 and P3_FIT301,
3581 we observe the following patterns: P2_FIT201, which is associated with
3582 incoming water flow and connected to valve P2_MV201, exhibits values
3583 greater than 2 when the valve is in the ON state. Conversely, when the
3584 valve is set to OFF, the sensor reading drops to 0.

3585 Regarding P3_FIT301, it appears to be linked to the behavior of the
3586 pump P3_P301 and correlates with the water flow out of the tank. When
3587 the pump is activated and in the ON state, the sensor records values greater
3588 than 2. On the other hand, when the pump is turned off and in the OFF
3589 state, the sensor reading returns to 0.

³⁵⁹⁰ **6.3.2.5 Properties**

³⁵⁹¹ Table 6.2 provides a summary of the properties inferred from the con-
³⁵⁹² jectures made throughout the different stages of the analysis.

#	Statement	Derived from
P23	The registers P3_DPIT301, P3_FIT301, and P3_LIT301 of PLC3 are considered measurements.	Preliminary Analysis Graphical Analysis
P24	The registers P3_MV301, P3_MV302, P3_MV303, P3_MV304, and P3_P302 od PLC3 are considered actuators.	Preliminary Analysis Graphical Analysis
P25	The register P3_P301 of PLC3 is considered a spare actuator.	Preliminary Analysis Graphical Analysis Invariant Analysis
P26	The register P3_LIT301 of PLC3 represents the level sensor of the tank controlled by PLC3.	Preliminary Analysis Graphical Analysis
P27	P2_MV201 and P3_P302 are 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 transitions 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 transitions 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 speed of thank 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 serves as 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

3593 6.3.3 Reverse Engineering of PLC3 and PLC4

3594 In the final phase of the reverse engineering process, the focus is di-
 3595 rected towards the subsystem consisting of PLC3 and PLC4 in the iTrust
 3596 SWaT system. Given the constraints of the thesis, this section will provide
 3597 a concise and schematic overview compared to the earlier sections.

3598 6.3.3.1 Pre-processing - Preliminary Analysis

3599 **Measurements and Actuators Recognition** Listing 6.11 shows the out-
 3600 comes obtained from automatic recognition of likely measurements and
 3601 actuators for PLC4. We omit those related to PLC3 as they are already
 3602 known.

```
3603 1 Actuators:  

  3604 2 ...  

  3605 3  

  3606 4 Sensors:  

  3607 5 ...  

  3608 6 P4_AIT401 {'max_lvl': 148.8, 'min_lvl': 148.8}  

  3609 7 P4_AIT402 {'max_lvl': 191.1, 'min_lvl': 185.5}  

  3610 8 P4_FIT401 {'max_lvl': 1.7, 'min_lvl': 1.7}  

  3611 9 P4_LIT401 {'max_lvl': 1002.8, 'min_lvl': 775.8}  

  3612 10
```

```

3613 11 Hardcoded setpoints or spare actuators:
3614 12 ...
3615 13 P4_P401 [1.0]
3616 14 P4_P402 [2.0]
3617 15 P4_P403 [1.0]
3618 16 P4_P404 [1.0]
3619 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 analogies from previous cases, P4_P402 is speculated to represent a pump that is constantly in the ON state and therefore active. However, regarding P4_UV401, it is unclear whether it is an actuator, a hardcoded setpoint, or a different type of register. Further analysis is needed to determine its exact purpose and functionality within the system.

On the other hand, it can be concluded that P4_P401, P4_P403, and

3644 P4_P404 are spare actuators, specifically pumps.

3645 **Actuator State Durations** Since there are no state-changing actuators within
3646 PLC4, further analysis regarding the duration of actuator states will not
3647 be performed for this subsystem. Please refer to Section 6.3.2.1 for evalua-
3648 tions of the duration of actuator states in PLC3.

3649 **Actuator State Changes** As previously mentioned, our assumption is
3650 that P4_LIT401 serves as the level sensor for the tank controlled by PLC4.
3651 In our analysis of PLC2-3, we speculated that P3_MV302 acted as a valve
3652 responsible for the incoming flow to an external element outside of that
3653 subsystem. To test this hypothesis, we examine the setpoints of P3_MV302
3654 in relation to the level of tank P4_LIT401. The setpoints of P3_MV302 corre-
3655 sponding to the tank level are presented in Listing 6.12.

```
3656 1 Actuator state changes:  
3657 2 ...  
3658 3 P3_MV302 prev_P3_MV302 P4_LIT401  
3659 4 0 2 1000.5510  
3660 5 0 1 784.4911  
3661 6 0 2 922.1866  
3662 7 0 1 881.0818  
3663 8 0 2 1000.2820  
3664 9 0 1 786.1061  
3665 10 0 2 922.8018  
3666 11 0 1 881.8508  
3667 12 ... ... ...
```

Listing 6.12: P3_MV302 state changes in relation to P4_LIT401

3668 The initial analysis suggests a potential correlation between the tank
3669 level values of P4_LIT401 and the behavior of P3_MV302. The ON state of
3670 P3_MV302 aligns with an increase in the tank level, while the OFF state cor-
3671 responds to a decrease. However, further analysis is required to provide
3672 additional evidence and support for this conjecture.

3673 **6.3.3.2 Graphs and Statistical Analysis**

3674 We will attempt to gain a deeper understanding of the pattern exhib-
 3675 ited by the PLC4 registers by referencing Figure 6.9.

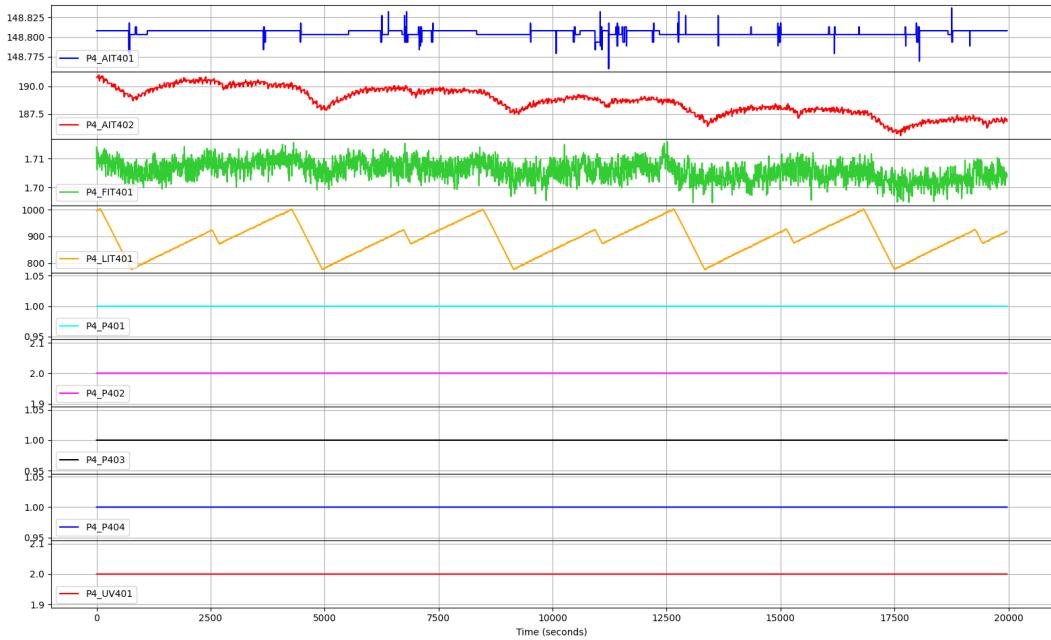


Figure 6.9: PLC4 registers

3676 The image reveals interesting behavior in the P4_AIT401 and P4_AIT401
 3677 registers. Notably, P4_AIT401 exhibits a linear trend rather than a cyclic
 3678 one, with values oscillating within a narrow range. This suggests that,
 3679 similar to P2_AIT201 (refer to Section 6.3.2.2), this register may correspond
 3680 to a sensor measuring a specific water property. On the other hand, P4_AIT402
 3681 appears to follow the level trend of sensor P4_LIT401, but with a down-
 3682 ward cyclic pattern where each cycle starts at a lower level than the pre-
 3683 vious one. Given the limited value range and the similarity in naming
 3684 conventions, it is highly likely that this register represents another water
 3685 property sensor rather than a tank.

3686 Additionally, P4_FIT401 does not display a cyclic pattern like the other
 3687 registers of the same type, and its values exhibit minimal variation, cor-
 3688 roborating the findings from the previous analysis phase.

3689 Now, let us refer to Figure 6.10 to seek confirmation regarding the con-
 3690 jecture that implicates valve P3_MV302 as the responsible actuator for tank
 3691 filling.

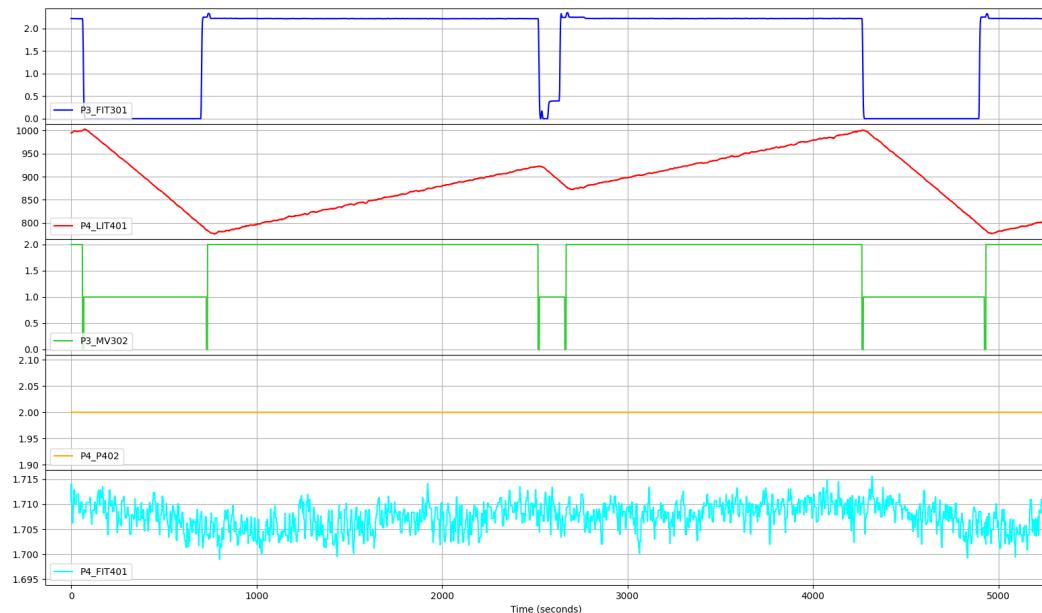


Figure 6.10: P3_MV302 and P4_LIT401 behaviors

3692 The image provides clear evidence that the behaviors of valve P3_MV302
 3693 and tank level sensor P4_LIT401 perfectly align. Additionally, P3_FIT301
 3694 (and its corresponding sensor, P3_DPIT301) appear to be related to the pat-
 3695 tern observed in P4_LIT401.

3696 Upon closer observation, it becomes apparent that the tank controlled
 3697 by PLC4 does **not have plateau periods**. When the incoming water flow
 3698 ceases, the tank immediately begins to empty. Based on our findings in
 3699 previous subsystems, we speculate that the actuator responsible for tank
 3700 emptying could be the pump indicated by register P4_P402. This specula-
 3701 tion is further supported by the nearly constant trend observed in sensor
 3702 P4_FIT401.

3703 Figure 6.11 depicts the correlation between the tanks within this sub-
 3704 system and the actuators that are responsible for their filling cycle.

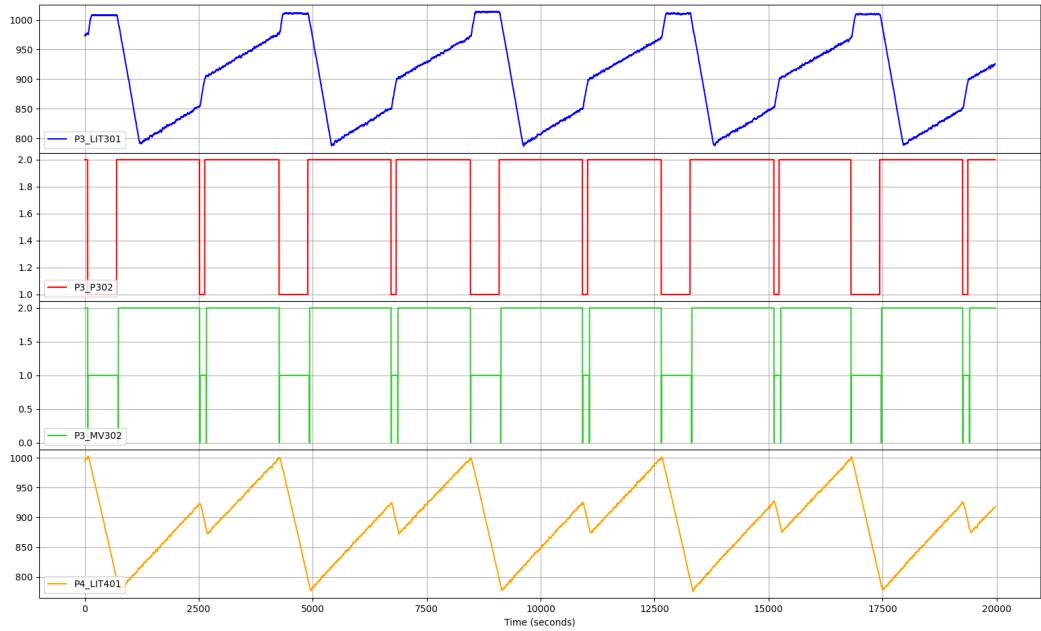


Figure 6.11: Tanks in subsystem PLC3-4 and their correlation.

3705 Further analysis was conducted to investigate whether valves P3_MV301,
 3706 P3_MV303, and P3_MV304 played a role in the tank filling cycle of PLC4.
 3707 However, the analysis did not confirm this hypothesis. Therefore, it can
 3708 be speculated that these valves are connected to other parts of the system
 3709 that are not currently discussed in this analysis.

3710 6.3.3.3 Invariant Inference and Analysis

3711 The invariant analysis for the subsystem consisting of PLC3-4 will be
 3712 brief as the states of the subsystem align with the states of valve P3_MV302,
 3713 with P4_P402 being constant throughout.

3714 **General Invariants** Again, the analysis of the general invariants offers
 3715 no new information compared to what we conjectured earlier. We there-
 3716 fore continue with the analysis of the current system configuration.

3717 **Analysis of the Current System Configuration** The analysis of the cur-
 3718 rent system configuration provides confirmation that when the system is

3719 in the state [$P3_MV302 == 1$, $P4_P402 == 2$], the tank level, as indicated
 3720 by sensor $P4_LIT401$, shows a decreasing slope. Additionally, it is noted
 3721 that the spare actuators align with the expected behavior.

3722 However, in the state [$P3_MV302 == 2$, $P4_P402 == 2$], the invariants
 3723 generated by Daikon do not provide the anticipated slope data, which was
 3724 expected to be increasing based on previous observations. This is due to
 3725 the descending slope between levels 880 and 925 of the tank represented
 3726 by $P4_LIT401$: by refining the analysis between the two increasing slope
 3727 sections we get the correct outcome, as shown in Listing 6.13.

```

3728 1 =====
3729 2 P3_MV302 == 2 && P4_P402 == 2 && P4_LIT401 < 990 &&
3730 3   ↪ P4_LIT401 > 930
3731 4 =====
3732 5 ...
3733 6 slope_P4_LIT401 == slope_P3_LIT301 == P4_P404 == P4_P403
3734 7   ↪ == P4_P401 == P3_P301 == P3_MV304 == P3_MV303 ==
3735 8   ↪ P3_MV301 == 1.0
3736 9 ...

```

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

3737 6.3.3.4 Business Process Mining and Analysis

3738 The process mining step for the PLC3-4 subsystem is straightforward.
 3739 In this phase, we will not focus on determining the chronological order of
 3740 states (as we already know they correspond to the states of valve $P3_MV302$).
 3741 Instead, we will examine the setpoints and extract any available informa-
 3742 tion concerning additional measurements. Figure 6.12 presents the activity
 3743 diagram depicting the subsystem under analysis.

3744 The diagram provides confirmation that when the system is in the state
 3745 [$P3_MV302 == 2$, $P4_P402 == 2$], the tank level trend is increasing. Con-
 3746 versely, when the system is in the state [$P3_MV302 == 2$, $P4_P402 == 2$],
 3747 the trend is decreasing.

3748 Regarding the setpoints, based on the standard deviation of the data

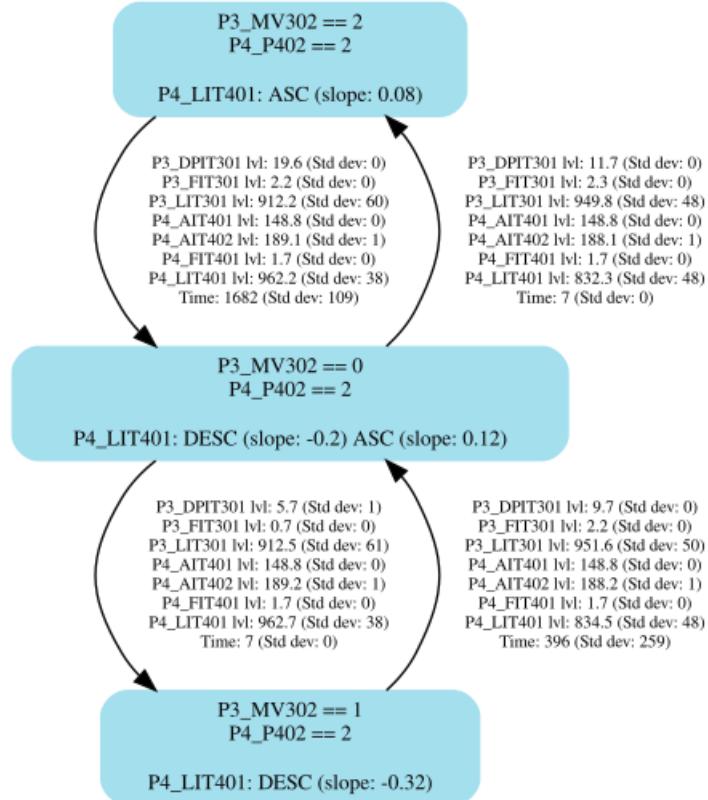


Figure 6.12: Activity diagram for PLC3-4

3749 from P4_LIT401, the absolute setpoints for the system are as follows: 785
 3750 (minimum), 880 (relative minimum), 925 (relative maximum), and 1000
 3751 (maximum).

3752 Furthermore, from the process mining phase, we can extract information
 3753 about P3_FIT301 and P3_DPIT301. P3_FIT301 registers values greater
 3754 than 2 when both pump P4_P402 and valve P3_MV302 are in the ON state,
 3755 while it drops to an average of 0.7 when the valve is off. Similarly, P3_DPIT301
 3756 registers values close to 20 when the valve is ON, and these values de-
 3757 crease when the valve is OFF.

3758 **6.3.3.5 Properties**

3759 Table 6.3 provides a summary of the properties inferred from the conjectures made throughout the different stages of the analysis. Regrettably, 3760 we encountered difficulties in determining the type and purpose of the 3761 register labeled as P4_UV401.

#	Statement	Derived from
P37	The registers P4_AIT401, P4_AIT402, P4_FIT401, and P4_LIT401 of PLC4 are considered measurements.	Preliminary Analysis Graphical Analysis
P38	The register P4_P402 of PLC4 is considered an actuator. Its status is constantly ON.	Preliminary Analysis Graphical Analysis Business Process
P39	The registers P4_P401, P4_P403, and P4_P404 of PLC4 are considered spare actuators.	Preliminary Analysis Graphical Analysis Invariant Analysis
P40	The register P4_LIT401 of PLC4 represents 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

Conclusion and Future Work

3763 **Discussion** At this stage, we have completed our (partial) reverse engi-
3764 neering of the iTrust SWaT System, aiming to achieve a sufficient level of
3765 the physical process, or *process comprehension*. This *process comprehension*
3766 enables us to plan a targeted attack on the system using the information
3767 obtained through the dynamic analysis conducted with the framework
3768 outlined in Chapter 4, employing a black box approach.

3769 To evaluate the accuracy of the information obtained about the SWaT
3770 system, we can refer to Figure 7.1 [66], which provides a schematic rep-
3771 resentation of the SWaT system. An x indicates placement of sensors.
3772 By comparing the information derived from our analysis with the system
3773 schematic, we can assess the validity and reliability of our findings.

3774 This image, while not comprehensive, serves to validate the accuracy
3775 of the properties derived from our system analysis of the first four stages
3776 of the SWaT system. It demonstrates that we have successfully identified
3777 the actuators and sensors within the system, and in some cases, we have
3778 even determined their specific roles within the physical process.

3779 Figure 7.2 [66] provides an alternative representation of the SWaT sys-
3780 tem from the perspective of the Human-Machine Interface (HMI). This
3781 depiction complements the previous diagram by adding more contextual
3782 information and enhancing overall understanding of the system. It offers

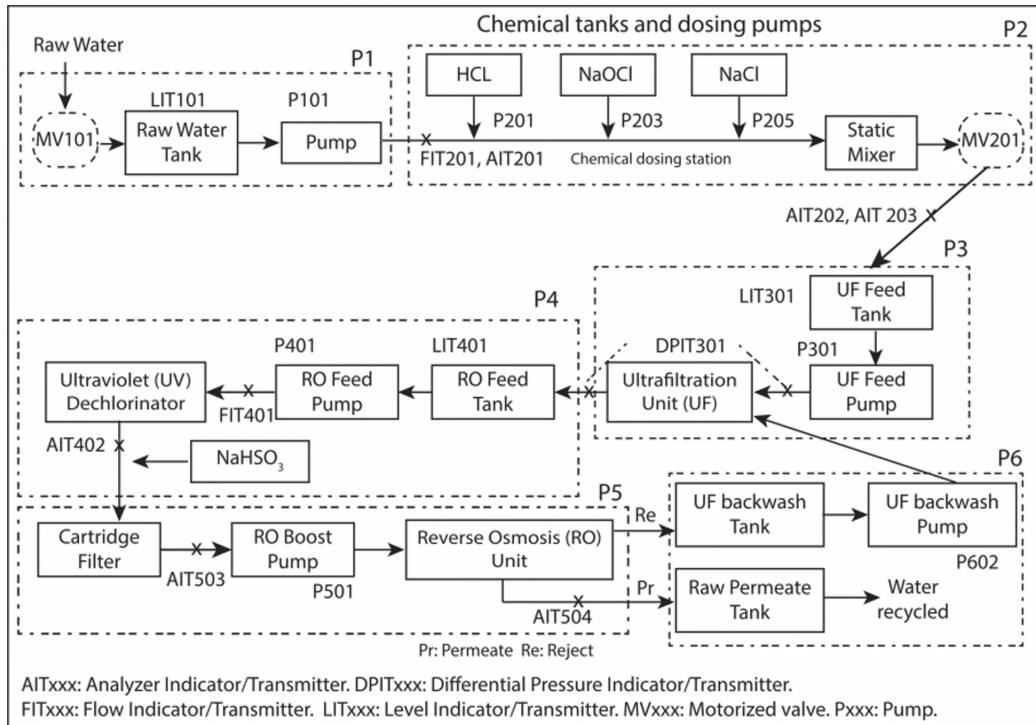


Figure 7.1: iTsut SWaT schema

3783 a comprehensive view of the system's components and their relationships,
 3784 thereby improving the clarity and comprehension of the system's structure
 3785 and functioning.

3786 This figure introduces additional elements that were not included in
 3787 the previous schematic shown in Figure 7.1. It incorporates spare actuators
 3788 and other components such as chemical tanks and the ultrafiltration
 3789 membrane, which were not explicitly represented as registers in the available
 3790 datasets. A comprehensive list of sensors and actuators, along with
 3791 their respective functions within the system, can be found in the paper by
 3792 Adepu et al. [69].

3793 However, this figure also serves as further confirmation of the effectiveness
 3794 of our analysis and the robustness of the framework we employed.
 3795 Thus, we can confidently assert that **we have successfully achieved the**
 3796 **objectives of this thesis**. We have not only validated the reverse engineering
 3797 methodology introduced by Ceccato et al., but also improved and

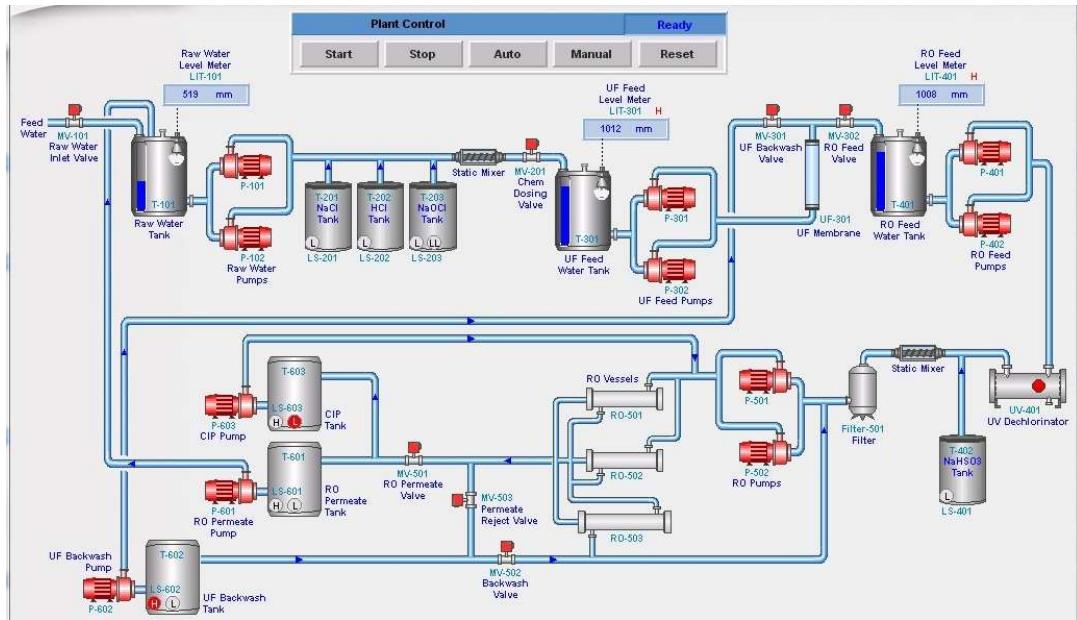


Figure 7.2: iTrust SWaT schema from HMI point of view.

enhanced it through its practical application, making it more adaptable and valuable in real-world contexts.

Future work The framework, along with the thesis sources and analysis files, is openly accessible on the dedicated GitHub repository [71]. There is potential for further improvement and expansion of the framework. One possibility is integrating the analysis framework with an additional framework that utilizes the gathered data to generate targeted system attacks, building upon the writer's previous work as described in Section 3.2.2.

Moreover, the proposed framework can benefit from various improvements and the introduction of specific new features to enhance its capabilities and effectiveness in analyzing industrial control systems. Here are some potential avenues for future work and enhancements to consider:

- enhance the scanning and data gathering phase by reimplementing it to include support for multiple industrial protocols, in addition to Modbus. This improvement would enable the framework to detect and communicate with PLCs using various protocols commonly

- 3814 used in ICS environments, expanding its compatibility and usability;
- 3815 • enhance the automatic recognition of sensors and actuators by lever-
- 3816 aging advanced machine learning and artificial intelligence techniques.
- 3817 This improvement would involve training models to accurately iden-
- 3818 tify and classify different types of sensors and actuators based on
- 3819 their data patterns and characteristics;
- 3820 • complete the implementation of network traffic data within the Busi-
- 3821 ness Process part of the framework. However, it is crucial to ensure
- 3822 that the network traffic data and physical process data used for anal-
- 3823 ysis are reliable and not compromised by external attackers, as accu-
- 3824 rate and untampered data is essential for effective analysis;
- 3825 • implement an automated system for recognizing real system config-
- 3826 urations by excluding actuators that do not significantly contribute
- 3827 to changing the system behavior within the analyzed PLC. This au-
- 3828 tomatic recognition system can save time and effort by eliminating
- 3829 the need to manually filter out non-contributing actuators during the
- 3830 analysis process.

List of Figures

2.1	ICS architecture schema	10
2.2	PLC architecture	13
2.3	PLC communication schema	14
2.4	Comparison between ST language and Ladder Logic	15
2.5	Example of HMI for a water treatment plant	18
2.6	Modbus Request/Response schema	20
2.7	Modbus RTU frame and Modbus TCP frame	21
2.8	Example of packet sniffing on the Modbus protocol	23
2.9	OSI model for EtherNet/IP stack	24
3.1	Workflow of Ceccato et al.'s stages and operations with used tools	32
3.2	The simplified SWaT system used for running Ceccato et al. methodology	33
3.3	Output graphs from graph analysis	37
3.4	An example of Disco generated activity diagram for PLC2 .	41
3.5	Execution traces of InputRegisters_IW0 on the three PLCs .	43
3.6	Business process with states and Modbus commands for the three PLCs	45
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.	51

3.8	Behavior of the Graph Analysis on the Ceccato et al.'s tool	53
3.9	Histogram plot overshadowing statistical information shown on the terminal window in the background	54
3.10	Example of Daikon's output	56
4.1	Network Communications between 192.168.1.40 (Source) and 192.168.1.30 (Destination)	72
4.2	Network Communications between 192.168.1.40 (Source) and 192.168.1.30 (Destination) in textual mode	73
4.3	Slope comparison with granularity 30 (a), 60 (b) and 120 seconds (c)	81
4.4	Savitzky-Golay filter (blue line), STL decomposition (green) and RDP algorithm (orange) comparison	83
4.5	Slope after the application of the STL decomposition	84
4.6	The new slope representation (green line) and the smoothed measurement data obtaind with the STL decomposition (red)	85
4.7	Plotting registers on the same y-axis	89
4.8	Example of the new graph analysis	90
4.8	Example of the new graph analysis (cont.)	91
4.9	Example of transitive closure graphs for invariants in PLC1 of iTrust SWaT system	95
4.10	Directory (a) and outcome files (b) for the single actuator states analysis	99
4.11	Daikon outcome files for system configuration analysis. Each file represents a single system state	99
4.12	Activity diagram for PLC1 of the iTrust SWaT system	106
5.1	SWaT architecture	113
5.2	Sensors and actuators associated with each PLC	114
5.3	SWaT network architecture and zoning	116
5.4	SWaT stabilization	118
6.1	Simplified graph of the iTrust SWaT system network	123
6.2	Chart of PLC1-2 registers	130

6.3	Chart of PLC1-2 registers without spare actuators (particular)	131
6.4	Slope calculation anomaly (in the circle)	136
6.5	Activity diagram for PLC1-2	137
6.6	Verifying the conjecture about valve P2_MV201	146
6.7	PLC3 registers	147
6.8	P3_MV301 and P3_MV303 analysis	148
6.9	PLC4 registers	158
6.10	P3_MV302 and P4_LIT401 behaviors	159
6.11	Tanks in subsystem PLC3-4 and their correlation.	160
6.12	Activity diagram for PLC3-4	162
7.1	iTrust SWaT schema	166
7.2	iTrust SWaT schema from HMI point of view.	167

List of Tables

2.1	differences between Information Technology (IT) and Industrial Control Systems (ICSs)	8
2.2	Modbus Function Codes list	22
3.1	Summary table of Ceccato et al. framework limitations	58
4.1	Summary table of proposed framework enhancements	110
5.1	main IP addresses of the six PLCs and SCADA in SWaT	116
5.2	SWaT network traffic data	119
6.1	Properties of the PLC1-2 subsystem	141
6.2	Properties of the PLC2-3 subsystem	155
6.3	Properties of the PLC3-4 subsystem	164

Listings

3.1	Example of registers capture	35
3.2	Example of raw network capture	36
3.3	Statistical data for PLC1_InputRegisters_IW0 register	38
3.4	Generic example of a .spinfo file for customizing rules in Daikon	39
3.5	The three sections of Daikon analysis outcomes	39
4.1	Novel Framework structure and Python scripts	63
4.2	Paths and parameters for the Pre-processing phase in <i>config.ini</i> file	75
4.3	config.ini parameters for dataset enriching	77
4.4	Example of preliminary system analysis	86
4.5	Standard Daikon output for PLC1 of the iTrust SWaT system	93
4.6	Revised Daikon output with transitive closures for PLC1 of the iTrust SWaT system	96
4.7	Daikon outcomes for the system configuration MV101 == 2, P101 == 1 on LIT101	100
4.8	Example of the data contained in the produced JSON file	104
6.1	Preliminary analysis outcomes for sensors and actuators of PLC1-2	126
6.2	Time duration of the states of actuators P1_MV101 and P1_MV201 of PLC1-2	128
6.3	P1_P101 state changes in relation to P1_LIT101	128
6.4	General Invariants for PLC1-2	133

6.5	Conditional Invariants for states 1 and 2 of P1_MV101	134
6.6	Preliminary analysis outcomes for sensors and actuators of PLC2-3	141
6.7	Time duration of the states of actuators of PLC3	142
6.8	P2_MV201 state changes in relation to P3_LIT301	144
6.9	Conditional Invariants for P2_MV201 and P3_P302	149
6.10	Slope calculation of P3_LIT301 for the 850-900 and 970-1000 intervals related to tank levels	152
6.11	Preliminary analysis outcomes for sensors and actuators of PLC3-4	155
6.12	P3_MV302 state changes in relation to P4_LIT401	157
6.13	Daikon manual analysis for P3_MV302 == 2	161

References

- [1] NIST. *ICS defintion*. URL: https://csrc.nist.gov/glossary/term/industrial_control_system (visited on 2023-04-06).
- [2] Wikipedia. *Stuxnet*. URL: <https://en.wikipedia.org/wiki/Stuxnet> (visited on 2023-05-18).
- [3] Wikipedia. *Industroyer*. URL: <https://en.wikipedia.org/wiki/Industroyer> (visited on 2023-05-18).
- [4] 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).
- [5] 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).
- [6] G. M. Makrakis, C. Kolias, 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>.
- [7] NIST. *PLC defintion*. URL: https://csrc.nist.gov/glossary/term/programmable_logic_controller (visited on 2023-04-06).

- [8] W. Bolton. *Programmable Logic Controllers, 6th edition*. Newnes, 2015, pp. 7–9.
- [9] 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.
- [10] 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/ICIInfs.2013.6731959>.
- [11] Wikipedia. *Remote Terminal Unit*. URL: https://en.wikipedia.org/wiki/Remote_terminal_unit (visited on 2023-04-08).
- [12] What is SCADA? URL: <https://www.inductiveautomation.com/resources/article/what-is-scada> (visited on 2023-04-05).
- [13] NIST. *HMI defintion*. URL: https://csrc.nist.gov/glossary/term/human_machine_interface (visited on 2023-04-11).
- [14] 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).
- [15] Schneider Electrics. *What is Modbus and how does it work?* URL: <https://www.se.com/us/en/faqs/FA168406/> (visited on 2023-04-13).
- [16] O. Morando. *Hacking: Modbus - Cyber security OT/ICS*. URL: <https://blog.omarmorando.com/hacking/ot-ics-hacking/hacking-modbus> (visited on 2023-04-15).
- [17] 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).

- [18] 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).
- [19] Odva, Inc. URL: <https://www.odva.org> (visited on 2023-04-21).
- [20] *Introduction to EtherNet/IP Technology*. URL: https://scadahacker.com/library/Documents/ICS_Protocols/Intro%20to%20EthernetIP%20Technology.pdf (visited on 2023-04-19).
- [21] ODVA Inc. *Common Industrial Protocol*. URL: <https://www.odva.org/technology-standards/key-technologies/common-industrial-protocol-cip/> (visited on 2022-12-26).
- [22] Wikipedia. *Common Industrial Protocol*. URL: https://en.wikipedia.org/wiki/Common_Industrial_Protocol (visited on 2023-04-21).
- [23] *What is the Common Industrial Protocol (CIP)?* URL: <https://www.motioncontroltips.com/what-is-the-common-industrial-protocol-cip/> (visited on 2023-04-21).
- [24] *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).
- [25] 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>.
- [26] 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>.
- [27] Wikipedia. *CODESYS*. URL: <https://en.wikipedia.org/wiki/CODESYS> (visited on 2023-05-18).

- [28] 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).
- [29] 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>.
- [30] Singapore University of Technology and Design. *iTrust - Center for Research in Cyber Security*. URL: <https://itrust.sutd.edu.sg/> (visited on 2023-05-15).
- [31] 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>.
- [32] *Association rules mining*. URL: https://en.wikipedia.org/wiki/Association_rule_learning (visited on 2023-04-21).
- [33] 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>.
- [34] 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).
- [35] Ceccato *et al.* *reverse engineering prototype tool*. URL: <https://github.com/SeUniVr/PLC-RE> (visited on 2023-04-23).

- [36] 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).
- [37] 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>.
- [38] M. Oliani. *AttackPLC - Project for Network Security teaching*. URL: <https://github.com/marcooliani/AttackPLC> (visited on 2023-04-23).
- [39] Ray - Productionizing and scaling Python ML workloads simply. URL: <https://www.ray.io/> (visited on 2023-04-25).
- [40] Tshark. URL: <https://tshark.dev/> (visited on 2023-04-25).
- [41] Wireshark. URL: <https://www.wireshark.org/> (visited on 2023-04-25).
- [42] pandas - Python Data Analysis library. URL: <https://pandas.pydata.org/> (visited on 2023-04-25).
- [43] NumPy - The fundamental package for scientific computing with Python. URL: <https://numpy.org/> (visited on 2023-04-25).
- [44] R project for statistical computing. URL: <https://www.r-project.org/> (visited on 2023-04-25).
- [45] SciPy - Fundamental algorithms for scientific computing with Python. URL: <https://scipy.org/> (visited on 2023-04-25).
- [46] The Daikon dynamic invariant detector. URL: <http://plse.cs.washington.edu/daikon/> (visited on 2023-04-25).
- [47] Enhancing Daikon output. URL: <https://plse.cs.washington.edu/daikon/download/doc/daikon/Enhancing-Daikon-output.html#Splitter-info-file-format> (visited on 2023-04-25).
- [48] Process mining. URL: https://en.wikipedia.org/wiki/Process_mining (visited on 2023-04-26).

- [49] *Fluxicon Disco*. URL: <https://fluxicon.com/disco/> (visited on 2023-04-26).
- [50] *OpenPLC - Open source PLC software*. URL: <https://openplcproject.com/> (visited on 2023-04-23).
- [51] *Docker*. URL: <https://www.docker.com/> (visited on 2023-04-26).
- [52] MathWorks. *Simulink*. URL: <https://it.mathworks.com/products/simulink.html> (visited on 2023-04-26).
- [53] 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).
- [54] *ProM Tools - The Process Mining Framework*. URL: <https://promtools.org/> (visited on 2023-04-26).
- [55] Fraunhofer Institute for Applied Information Technology. *pm4py - Process Mining for Python*. URL: <https://promtools.org/> (visited on 2023-04-26).
- [56] statsmodels. *statsmodels - statistical models, hypothesis tests, and data exploration*. URL: <https://www.statsmodels.org/stable/index.html> (visited on 2023-05-05).
- [57] NetworkX. *NetworkX - Network Analysis in Python*. URL: <https://networkx.org/> (visited on 2023-05-05).
- [58] Wikipedia. *Polynomial regression*. URL: https://en.wikipedia.org/wiki/Polynomial_regression (visited on 2023-05-21).
- [59] Wikipedia. *Decomposition of time series*. URL: https://en.wikipedia.org/wiki/Decomposition_of_time_series (visited on 2023-05-21).
- [60] M. Fleischmann. "Line simplification algorithms". In: (2020-04-27). URL: <https://martinfleischmann.net/line-simplification-algorithms/> (visited on 2023-05-21).
- [61] Wikipedia. *Savitzky-Golay filter*. URL: https://en.wikipedia.org/wiki/Savitzky-Golay_filter (visited on 2023-05-06).

- [62] Otext. *STL decomposition*. URL: <https://otexts.com/fpp2/stl.html> (visited on 2023-05-06).
- [63] Wikipedia. *Ramer-Douglas-Peucker algorithm*. URL: https://en.wikipedia.org/wiki/Ramer-Douglas-Peucker_algorithm (visited on 2023-05-21).
- [64] *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).
- [65] Wikipedia. *Transitive closures*. URL: https://en.wikipedia.org/wiki/Transitive_closure (visited on 2023-05-11).
- [66] 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).
- [67] 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>.
- [68] iTrust Singapore University of Technology and Design. *Datasets*. URL: https://itrust.sutd.edu.sg/itrust-labs_datasets/ (visited on 2023-05-15).
- [69] 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.
- [70] 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).
- [71] M. Oliani. *Master Degree's Thesis Framework and Sources*. URL: <https://github.com/marcooliani/Thesis> (visited on 2023-06-15).