

UNIVERSITY OF VERONA

Department of COMPUTER SCIENCE

Master's Degree in
COMPUTER SCIENCE AND ENGINEERING

Master Thesis

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

Candidate:

Marco OLIANI
VR457249

Supervisor:

Prof. Massimo MERRO

Co-supervisor:

Prof. Ruggero LANOTTE
University of Insubria

Academic Year 2022/2023

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

Abstract

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

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

This thesis analyzes Ceccato et al.'s reverse engineering-based methodology for ICS analysis [1], identifying limitations in its application to real-world scenarios. The framework's flexibility, network analysis, pre-processing, graphical analysis, invariant analysis, and business process analysis are improved to address these limitations and new functionalities are added. The enhanced framework is tested on a complex case study, the iTrust SWaT system, a scaled-down water treatment plant.

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

Contents

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

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

5 Case study: the iTrust SWaT System	113
5.1 Architecture	113
5.1.1 Physical Process	114
5.1.2 Control and Communication Network	116
5.2 Datasets	118
5.2.1 Our Case Study: the 2015 Dataset	119
6 Our Framework at Work on the iTrust SWaT System	123
6.1 Preliminary Operations	124
6.2 Planning the Analysis Strategy	124
6.3 Reverse Engineering of the iTrust SWaT System	126
6.3.1 Reverse Engineering of PLC1 and PLC2	128
6.3.1.1 Pre-processing - Preliminary Analysis . . .	128
6.3.1.2 Graphs and Statistical Analysis	132
6.3.1.3 Invariant Inference and Analysis	135
6.3.1.4 Business Process Mining and Analysis . .	138
6.3.1.5 Properties	141
6.3.2 Reverse Engineering of PLC2 and PLC3	143
6.3.2.1 Pre-processing - Preliminary Analysis . . .	143
6.3.2.2 Graphs and Statistical Analysis	147
6.3.2.3 Invariant Inference and Analysis	151
6.3.2.4 Business Process Mining and Analysis . .	154
6.3.2.5 Properties	157
6.3.3 Reverse Engineering of PLC3 and PLC4	158
6.3.3.1 Pre-processing - Preliminary Analysis . . .	158
6.3.3.2 Graphs and Statistical Analysis	161
6.3.3.3 Invariant Inference and Analysis	163
6.3.3.4 Business Process Mining and Analysis . .	164
6.3.3.5 Properties	166
6.4 PLCs Architecture	167
7 Conclusion and Future Work	169
Appendix A - Framework Scripts Handbook	175

List of Figures	179
List of Tables	183
Listings	185
References	187

Introduction

THE advent of Industry 4.0 has sparked significant transformations in the realm of *Industrial Control Systems* (ICS), physical and engineered systems whose operations are monitored, coordinated, controlled, and integrated by a computing and communication core [2]. They are used to control industrial processes such as manufacturing, product handling, production, and distribution [3]. In recent years, there has been a rapid expansion in the interconnectivity and convergence of IT (*Information Technology*) and OT (*Operational Technology*) systems. While this integration offers undeniable benefits in terms of operational efficiency and effectiveness of ICSs, it has also exposed these systems to the prevalent vulnerabilities commonly associated with IT environments. Consequently, there has been a parallel rise in **cyber-physical attacks**, which originates in the cyber environment and target the physical processes of these systems. These attacks aim to manipulate or disrupt the normal operation of ICSs, posing a considerable risk to their integrity and functionality.

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

monitoring, control, and maintenance purposes [4].

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

OT networks encompass two primary sub-networks: the *supervisory control network*, which connects the PLCs and HMIs, and the *field communications network*, which establishes connections between the PLCs and the associated field devices.

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

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

led to a transition from serial-based communication protocols to IP-based protocols. These protocols, known as **industrial protocols**, have become prevalent in modern industrial systems. Examples of industrial protocols include Modbus [5], DNP3 [6], EtherNet/IP [7], OPC UA [8], and S7comm [9] (which is a proprietary and closed protocol).

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

Open Problems To execute a successful cyber-physical attack, an attacker must possess a comprehensive understanding of the target system's characteristics and behavior, commonly referred to as *process comprehension* [10]. This entails possessing knowledge about various aspects, including the **operational domain** (such as water distribution or power generation), the controllers used (such as PLCs, RTUs, etc.), the **network topology** associated with them, relevant **measurements** within the facility (such as pressure, temperature, etc.), as well as the exposed physical components like **sensors and actuators** that can be vulnerable to attacks. By attaining such insights, the attacker can effectively identify vulnerabilities, exploit weaknesses, and manipulate the system in a targeted and *stealthy* manner. Moreover, to bypass certain types of *intrusion detection systems* (IDSs), the attacker needs to possess a general understanding of the **physical invariants** of the system.

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

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

- **System Architecture:** attackers aim to understand the architecture of the ICS, including the components involved, such as sensors, actuators, controllers, and human-machine interfaces. This knowledge helps them identify potential entry points and vulnerabilities within the system;
- **Industrial Protocols:** familiarity with industrial communication protocols, such as Modbus [5], DNP3 [6], or EtherNet/IP - CIP [11], is crucial for attackers. Understanding the protocols enables them to analyze the communication between different ICS components and potentially exploit vulnerabilities or manipulate the data flow;
- **Operational Processes:** attackers seek to comprehend the operational processes and workflows within the targeted ICS. This includes understanding the sequence of actions, system states, and interactions between various components. Such knowledge enables attackers to identify critical points where disruptions or malicious actions can have significant consequences;
- **Process Dependencies:** attackers analyze the dependencies between different physical processes or systems within the ICS. They aim to identify dependencies that, if disrupted or compromised, could have a cascading effect on the overall operation. This helps them strategize targeted attacks for maximum impact;
- **Attack Surface:** by evaluating the attack surface of the ICS, including network connectivity, remote access options, and third-party integrations, attackers can identify potential entry points and avenues for exploitation;
- **Attack Goals:** attackers define their specific goals, whether it be disruption of operations, theft of sensitive data, sabotage, or any other malicious intent. Understanding the desired outcomes helps them

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

Attackers can acquire a good level of process comprehension of an ICS through several methods. For instance, they can utilize **probing** techniques [12], which involve introducing slight perturbations to the system and observing its response and subsequent return to its original state. By analyzing these reactions, attackers can gain valuable insights into the system's behavior.

Another technique is the ***Man-in-the-Middle* approach** [10], where attackers intercept and analyze network communications to gather information about the system. This allows them to conduct **reconnaissance** of the system while avoiding detection.

Reverse engineering [1][13] is another method attackers can employ. By scanning memory logs of network-exposed PLCs and analyzing network traffic related to the industrial protocols, they can approximate the physical system's model. This technique enables them to derive useful information without directly intervening in the physical system, as they can work offline using collected logs and analysis tools.

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

Contribution

The primary objective of this thesis is to conduct an analysis of Ceccato et al.'s reverse engineering-based methodology [1]. The aim is to examine its strengths and limitations and apply it to a real-world case study.

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

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

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

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

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

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

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

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

The *Graphical Analysis* component has undergone significant improvements to enhance its effectiveness. **Subplots** now enable the visualization of registers as a whole, facilitating the identification of relationships between them. Users can select specific registers to view and zoom in on specific areas of the graphs without losing information.

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

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

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

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

Outline

The thesis is structured as follows:

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

Chapter 3: after an introductory section that provides a brief overview of the existing literature on process comprehension in industrial control systems, this chapter focuses on a specific paper that outlines a methodology to attaining process comprehension of an industrial system by employing dynamic blackbox analysis;

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

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

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

Chapter 7: outlines final conclusions and future work.

Background on Industrial Control Systems

INDUSTRIAL Control Systems (ICSs) are physical and engineered systems whose operations are monitored, coordinated, controlled and integrated by a computing and communication core [2]. They are used to control industrial processes such as manufacturing, product handling, production, and distribution [3]. ICSs are often found in critical infrastructure facilities such as power plants, oil and gas refineries, and chemical plant.

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

ICSs are becoming increasingly connected to the internet and other networks, which has led to increased concerns about their security. Industrial systems were not originally designed with security in mind, and many of them have known vulnerabilities that could be exploited by attackers. Additionally, the use of legacy systems and equipment can make it difficult to implement security measures. As a result, ICSs are increasingly seen

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

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

Table 2.1 summarizes the differences between traditional IT and ICSs [4]:

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

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

2.1 Industrial Control Systems Architecture

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

General ICS architecture consists in **six levels** each representing a functionality: this architecture comprising the OT and IT parts is represented in Figure 2.1 [19][4], according to the *Purdue Enterprise Reference Architecture* (PERA), or simply **Purdue Model**:

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

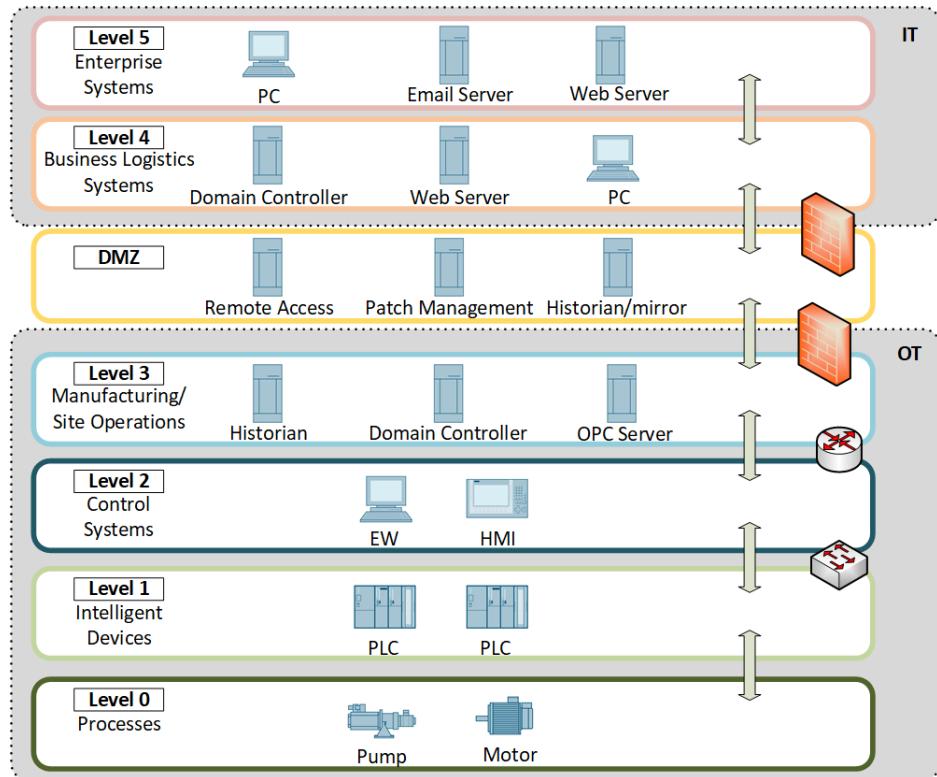


Figure 2.1: ICS architecture schema

- **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.
- 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.
- Level 5 (**Enterprise Systems**): represents the enterprise network, used for the business-to-business activities and for business-to-client pur-

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.

2.2.2 Controller Network Layer

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

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

2.2.2.1 Programmable Logic Controllers

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

Compared to relay systems and personal computers, PLCs are optimized for control tasks and industrial environments: they are rugged and designed to withdraw harsh conditions such as dust, vibrations, humidity and temperature: they have more reliability than personal computers, which are more prone to crash, and they are more compact and require less maintenance than a relay system.

Furthermore, I/O interfaces are already on the controller, so PLCs are easier to expand with additional I/O modules (if in a rack format) to manage more inputs and outputs, without reconfiguring hardware as in relay systems when a reconfiguration occurs.

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

PLC Architecture The basic hardware architecture of a PLC consists of these elements [21]:

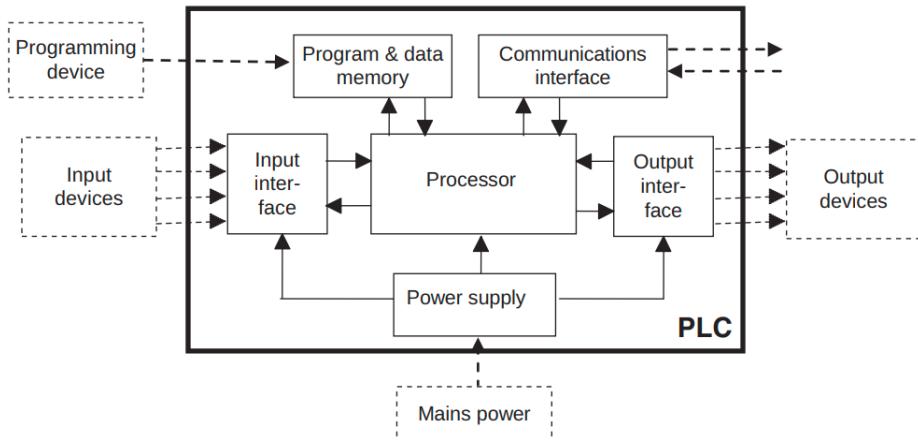


Figure 2.2: PLC architecture

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

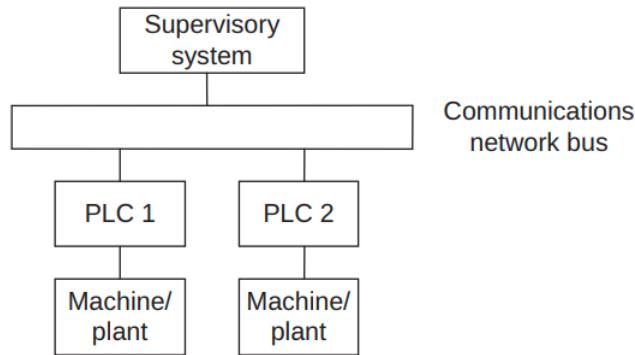


Figure 2.3: PLC communication schema

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

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

The user program needs to be uploaded on the PLC via the programming device and runs on the process image table in *scan cycles*: each scan is made up of three phases [1]:

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

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

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

```

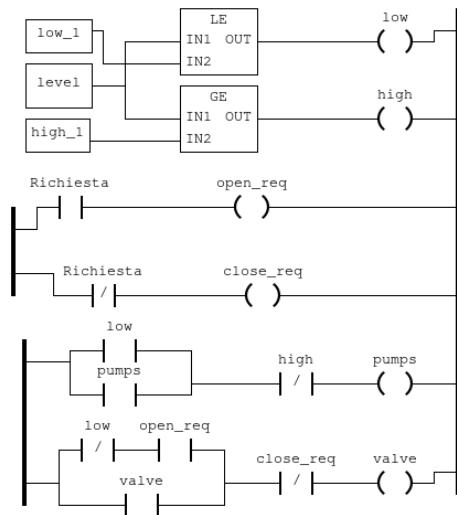
PROGRAM PLC1
VAR
    level AT %IW0 : INT;
    Richiesta AT %QX0..2 : BOOL;
    request AT %IW1 : INT;
    pumps AT %QX0..0 : BOOL;
    valve AT %QX0..1 : BOOL;
    low AT %MW0..0 : BOOL;
    high AT %MW0..1 : BOOL;
    open_req AT %MW0..3 : BOOL;
    close_req AT %MW0..4 : BOOL;
    low_1 AT %MW0 : INT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

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

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

```

(a) Example of ST programming



(b) Example of Ladder Logic

Figure 2.4: Comparison between ST language and Ladder Logic

PLC Security PLCs were originally designed to operate as closed systems, not connected and exposed to the outside world via communication networks: the question of the safety of these systems, therefore, was not a primary aspect. The advent of Internet has brought undoubted advantages, but has introduced problems relating to the safety and protection of PLCs from external attacks and vulnerabilities.

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

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

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

In addition to this, keeping the process network and Internet separated, limiting the use of USB devices among users to reduce the risks of infections, and using strong account management and maintenance policies are best practices to prevent attacks and threats and to avoid potential damages.

2.2.2.2 Remote Terminal Units

Remote Terminal Units (RTUs) are computers with radio interfacing similar to PLCs: they transmit telemetry data to the control center or to the PLCs and use messages from the master supervisory system to control connected objects [23].

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

Industries that require RTUs often operate in areas without reliable access to the power grid or require monitoring and control substations in remote locations. These include telecommunications, railways, and utilities that manage critical infrastructure such as power grids, pipelines, and water treatment facilities. The advanced technology of RTUs allows these industries to maintain essential services, even in challenging environments or under adverse weather conditions.

2.2.3 Area Control Layer

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

2.2.3.1 Supervisory Control And Data Acquisition

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

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

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

2.2.3.2 Human-Machine Interface

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

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

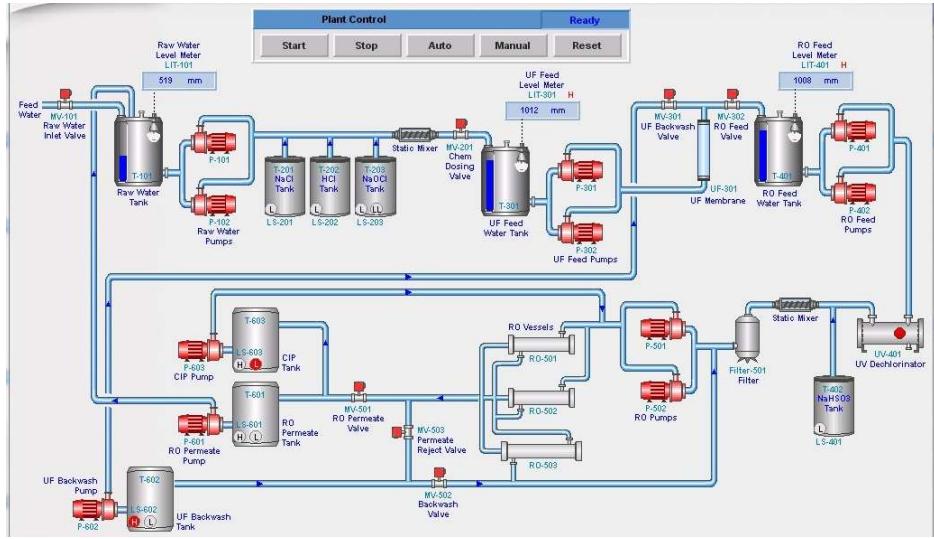


Figure 2.5: Example of HMI for a water treatment plant

2.2.4 Operations/Control Layer

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

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

2.2.5 Demilitarized Zone

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

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

2.2.6 Industrial Protocols

Industrial Protocols are the networks that are used to connect the different components of the ICS and allow them to communicate with each other. Industrial Protocols can include wired and wireless networks, such as Ethernet/IP, Modbus, DNP3, Profinet and others.

As mentioned at the beginning of this Chapter, industrial systems differ from classical IT systems in the purpose for which they are designed: controlling physical processes the former, processing and storing data the latter. For this reason, ICSs require different communication protocols than traditional IT systems for real time communications and data transfer.

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

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

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

2.2.6.1 Modbus

Modbus is a serial communication protocol developed by Modicon (now Schneider Electric) in 1979 for use with its PLCs [5] and designed expressly for industrial use: it facilitates interoperability of different devices connected to the same network (sensors, PLCs, HMIs, ...) and it is also often used to connect RTUs to SCADA acquisition systems.

Modbus is the most widely used communication protocol among industrial systems because it has several advantages:

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

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

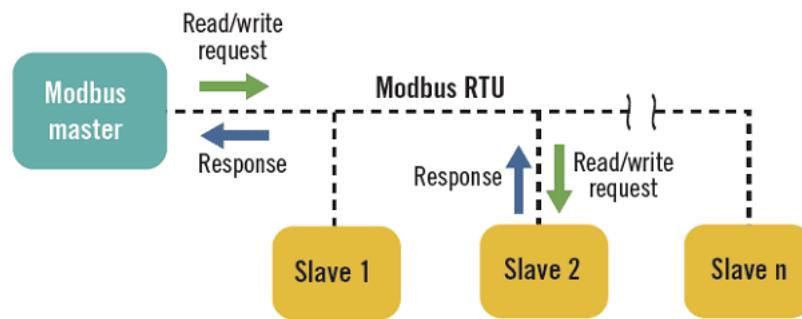


Figure 2.6: Modbus Request/Response schema

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

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

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

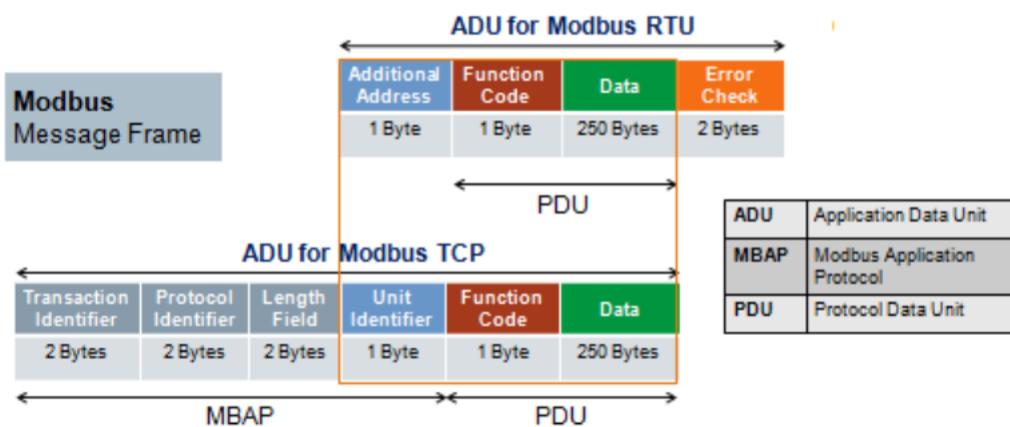


Figure 2.7: Modbus RTU frame and Modbus TCP frame

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

- *Coil*: binary type, read/write accessible by both masters and slaves
- *Discrete Input*: binary type, accessible in read-only mode by masters and in read/write mode by slaves
- *Analog Input*: 16 bits in size (word), are accessible in read-only mode by masters and in read/write mode by slaves
- *Holding Register*: 16 bits in size (word), accessible in read/write mode by both masters and slaves. Holding Registers are the most commonly used registers for output and as general memory registers.

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

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

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

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

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

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

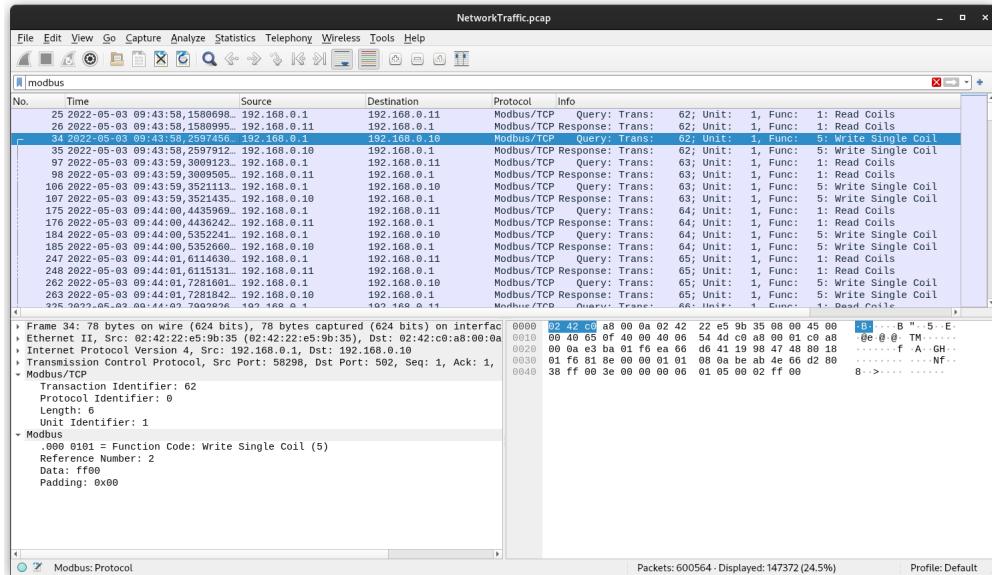


Figure 2.8: Example of packet sniffing on the Modbus protocol

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

2.2.6.2 EtherNet/IP

EtherNet/IP (where IP stands for *Industrial Protocol*) is an open industrial protocol that allows the *Common Industrial Protocol* (CIP) to run on a typical Ethernet network [11]. It is supported by ODVA [29].

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

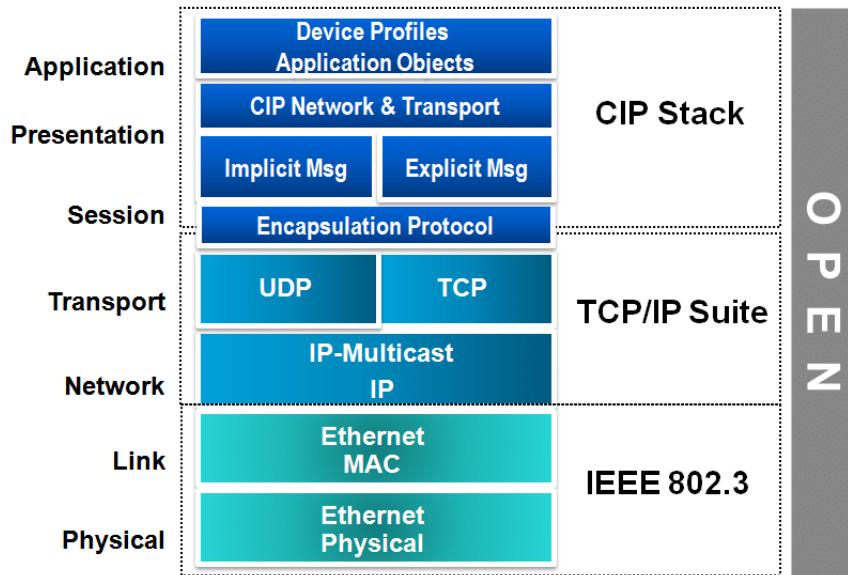


Figure 2.9: OSI model for EtherNet/IP stack

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

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

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

Transport layer At the transport level, EtherNet/IP encapsulates messages from the CIP stack into an Ethernet message, so that messages can be transmitted from one node to another on the network using the TCP/IP protocol. EtherNet/IP uses two forms of messaging, as defined by CIP standard [11][7]:

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

nected messaging uses TCP/IP to transmit messages across the network asking for connection resource each time from the *Unconnected Message Manager* (UCMM).

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

EtherNet/IP has two types of message connection [11]:

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

Session, Presentation and Application layer At the upper layers, EtherNet/IP implements the CIP protocol stack. We will discuss this protocol more in detail in Section 2.2.6.3.

2.2.6.3 Common Industrial Protocol (CIP)

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

CIP has been adapted in different types of network:

- **EtherNet/IP**, adaptation to *Transmission Control Protocol* (TCP) technologies
- **ControlNet**, adaptation to *Concurrent Time Domain Multiple Access* (CTDMA) technologies
- **DeviceNet**, adaptation to *Controller Area Network* (CAN) technologies
- **CompoNet**, adaptation to *Time Division Multiple Access* (TDMA) technologies

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

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

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

State of the Art

IN conventional IT systems, the objective of an attacker is to comprehend 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).

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.

In the context of OT systems, the notion of *reverse engineering* is not limited to its conventional definition, but also includes the concept of **process comprehension**. This term, introduced by Green et al. [10], refers to gaining a comprehensive understanding of the underlying physical process.

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

3.1 Literature on Process Comprehension

Keliris and Maniatikos The first approach presented in this section is by Keliris and Maniatikos [13]: they present a methodology for automating the reverse engineering of ICS binaries based on a *modular framework* (called ICSREF) that can reverse binaries compiled with CODESYS [34], one of the most popular and widely used PLC compilers, irrespective of the language used.

Yuan et al. Yuan et al. [35] propose a *data-driven* approach to discovering cyber-physical systems process behavior from data directly: to achieve this goal, they have implemented a framework whose purpose is to identify physical systems and transition logic inference, and to seek to understand the mechanisms underlying these processes, making furthermore predictions concerning their state trajectories based on the discovered models.

Feng et al. Feng et al. [36] developed a framework that can generate system *invariant rules* based on machine learning and data mining techniques from ICS operational data log. These invariants are then selected by systems engineers to derive IDS systems from them.

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

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

System as testbed.

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

Winnicki et al. Winnicki et al. [12] instead propose a different approach to process comprehension based on the *attacker's perspective* and not limited to mere *Denial of Service* (DoS): their approach is to discover the dynamic behavior of the system, in a semi-automated and process-aware way, through *probing*, that is, slightly perturbing the cyber physical system and observing how it reacts to changes and how it returns to its original state. The difficulty and challenge for the attacker is to perturb the system in such a way as to achieve an observable change, but at the same time avoid this change being seen as a system anomaly by the IDSs.

Green et al. Green et al. [10] also adopt an approach based on the attacker's perspective: this approach consists of two practical examples in a *Man in the Middle* (MitM) scenario to obtain, correlate, and understand all the types of information an attacker might need to plan an attack to alter the process while avoiding detection.

The paper shows *step-by-step* how to perform a ICS **reconnaissance**, a phase specifically designed to gather extensive intelligence on multiple fronts, including human factors, network and protocol information, details about the manufacturing process, industrial applications, and potential vulnerabilities. The primary goal is to accumulate a wealth of information to enhance understanding and awareness in these areas [40]).

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

Ceccato et al. Ceccato et al. [1] propose a methodology based on a *black box dynamic analysis* of an ICS using a reverse engineering tool to

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

derive from the scans performed on the memory registers of the exposed PLCs and the Modbus protocol network scans an approximate model of the physical process. This model is obtained by inferring statistical properties, business process and system invariants from data logs.

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

In the next section we will examine this latest work in more detail, which will be the basis for my work and thus the subsequent chapters of this thesis.

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

As previously mentioned, the paper introduces a methodology that relies on black box dynamic analysis of an Industrial Control System (ICS) and more particularly of its OT network. This methodology involves identifying potential Programmable Logic Controllers (PLCs) within the network and scanning the memory registers of these identified controllers. The purpose of this process is to obtain an approximate model of the controlled physical process.

The primary goal of this black box analysis is to establish a correlation between the different memory registers of the targeted PLCs and fundamental concepts of an OT network such as sensor values (i.e., measurements), actuator commands, setpoints (i.e., range of values of a physical variable), network communications, among others.

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

The second goal is to establish a relationship between the dynamic evolution of these fundamental concepts.

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

1. **scanning of the system and data pre-processing:** this phase involves gathering data to generate data logs for the registers of PLCs and for Modbus network communications.
2. **graphs and statistical analysis:** The collected data is utilized to provide insights into the memory registers associated with the Modbus protocol by leveraging graphs and statistical data. This analysis approach offers valuable information about the characteristics and patterns of the memory registers.
3. **invariants inference and analysis:** generates system invariants, which are used to identify specific patterns and regularities within the system. Additionally, this phase provides users with the capability to view invariants related to a particular sensor or actuator.
4. **business process mining and analysis:** Using event logs, this phase involves reconstructing the business process that depicts how a process is executed. This step enables a thorough understanding of the sequence of events that occur in the system and how they are interrelated, ultimately leading to a comprehensive overview of the business process.

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

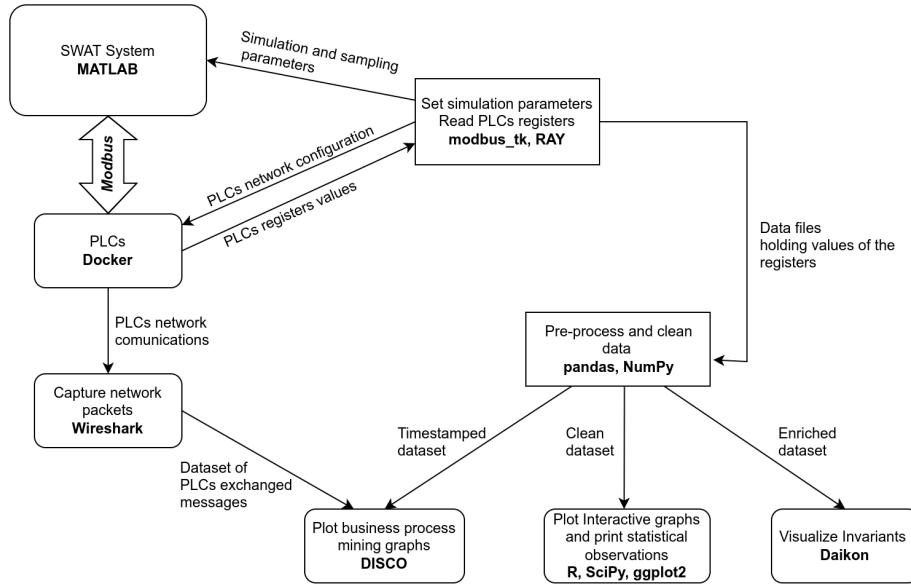


Figure 3.1: Workflow of Ceccato et al.'s stages and operations with used tools

3.2.1 Testbed

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

Stage 1 During the initial stage, a **tank** referred to as T-201 with a capacity 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 accumulated water from the tank, directing it to the next stage. Initially, the water flows from the T-201 tank to the *filtration unit* (which is not specifically identified by any sensor), and subsequently to a **second tank** denoted as T-202, with a capacity of 20 gallons.

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

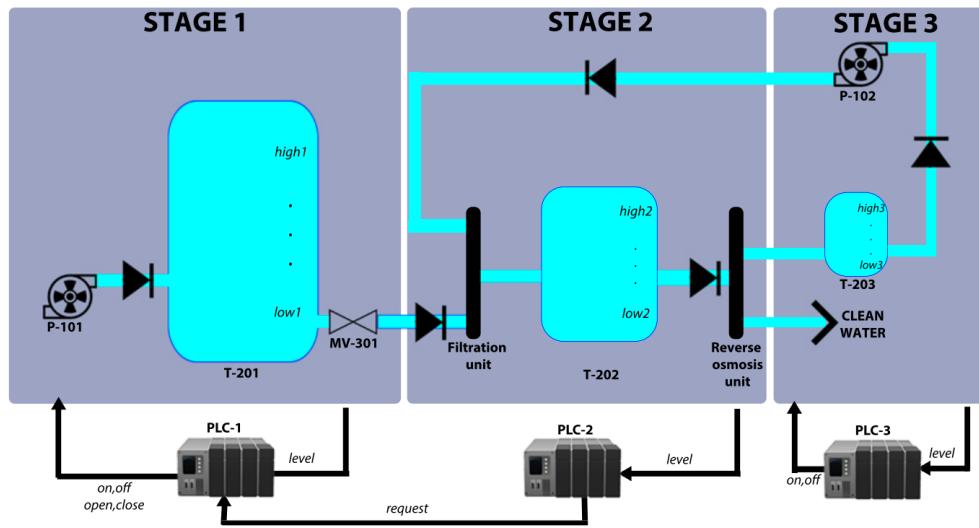


Figure 3.2: The simplified SWaT system used for running Ceccato et al. methodology

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

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

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

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

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

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

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

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

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

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

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

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

3.2.2 Scanning of the System and Data Pre-processing

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

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

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

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

Listing 3.1: Example of registers capture

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

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

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

Listing 3.2: Example of raw network capture

Data Pre-processing The data collected by scanning the memory registers of the PLCs are then reprocessed by a Python script and converted in order to create a distinct raw dataset in *Comma Separated Value* format (CSV) for each PLC, containing the memory register values associated with the corresponding controller registers. These datasets are reprocessed again through the Python modules for **pandas** [47] and **NumPy** [48] by another script to first perform a **data cleanup**, removing all unused registers, **merged** into a single dataset, and finally **enriched** with additional data, such as the **previous value** of all registers and the **measurement slope**, that is, the trend of the water level in the system tanks along the 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.

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

3.2.3 Graphs and Statistical Analysis

The paper mentions the presence of a *mild graph analysis*, performed using the framework **R** [49] for statistical analysis at the time of data gathering to find any uncovered patterns, trends and identify measurements and/or actuator commands through the analysis of registers holding mutable values.

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

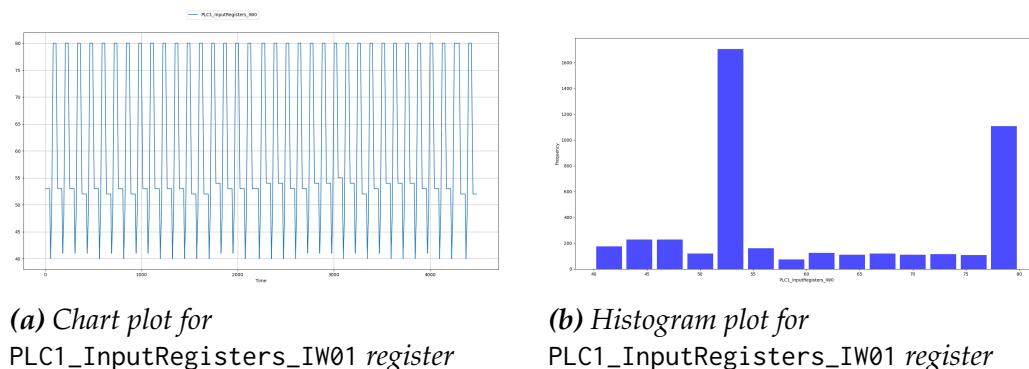


Figure 3.3: Output graphs from graph analysis

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

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

- the mean, median, standard deviation, maximum value and minimum value
- two tests for the statistical distribution: *Chi-squared* test for uniformity and *Shapiro-Wilk* test for normality, as shown in Listing 3.3:

```
Chi-squared test for uniformity
Distance      pvalue      Uniform?
12488.340    0.00000000    NO

Shapiro-Wilk test for normality
Test statistic   pvalue      Normal?
0.844        0.00000000    NO

Stats of PLC1_InputRegisters_IW0
Sample mean = 60.8881; Stddev = 13.0164; max = 80; min =
→ 40 for 4488 values
```

Listing 3.3: Statistical data for PLC1_InputRegisters_IW0 register

3.2.4 Invariant Inference and Analysis

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

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

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

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

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

Therefore, Daikon is fed with the enriched dataset obtained in the pre-processing phase²: a simple bash script launches Daikon (optionally specifying the desired condition for analysis in the *.spininfo* file), which output is simply redirected to a text file containing the general invariants of the system (i.e., valid regardless of any custom condition specified), those generated based on the custom condition in the *.spininfo* file, and those generated based on the negation of the condition (see Listing 3.5 below).

```
=====
aprogam.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 }
[...]
=====
aprogam.point:::POINT; condition="PLC1_InputRegisters_IW0
    ↪ > 60"
PLC1_InputRegisters_IW0 > PLC1_MemoryRegisters_MW0
PLC1_InputRegisters_IW0 > PLC1_Min_safety
PLC1_MemoryRegisters_MW0 < prev_PLC1_InputRegisters_IW0
[...]
=====
aprogam.point:::POINT; condition="not(
    ↪ PLC1_InputRegisters_IW0 > 60)"
PLC1_InputRegisters_IW0 < PLC1_MemoryRegisters_MW1
PLC1_InputRegisters_IW0 < PLC1_Max_safety
```

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

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

```
PLC1_MemoryRegisters_MW1 > prev_PLC1_InputRegisters_IW0  
[...]
```

Listing 3.5: The three sections of Daikon analysis outcomes

When the analysis is finished, the user is asked to enter the name of a registry to view its related invariants.

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

- measurements bounded by some setpoint;
- actuators state changes occurred in the proximity of setpoints or, vice versa, proximity of setpoints upon the occurrence of an actuator state change;
- state invariants of some actuators correspond to a specific trend in the evolution of the measurements (ascending, descending, or stable) or, vice versa, the measurements trend corresponds to a specific state invariant of some actuators.

3.2.5 Business Process Mining and Analysis

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

In the considered system, process mining begins by analyzing the event logs derived from scanning the memory registers of the PLCs and monitoring the network communications associated with the Modbus protocol, as detailed in Subsection 3.2.2. These event logs serve as the *execution trace* of the system. A Java program is utilized to extract and consolidate information from these event logs, resulting in a CSV format file that captures the relevant data.

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

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

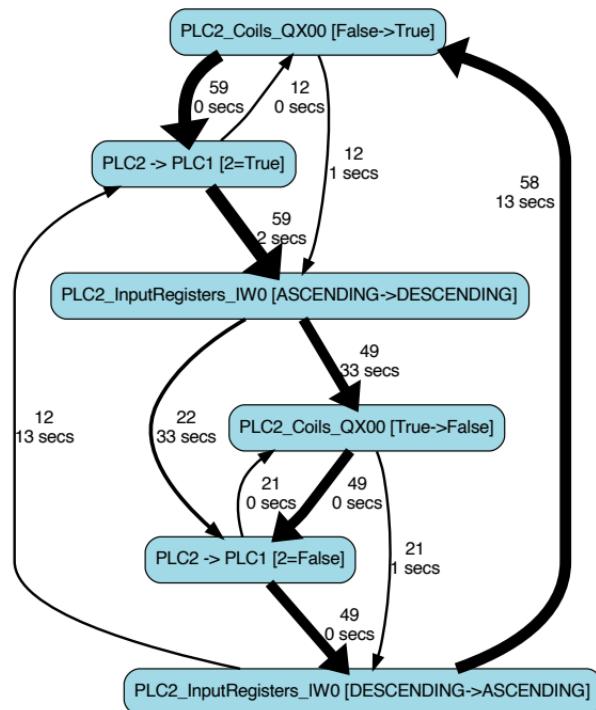


Figure 3.4: An example of Disco generated activity diagram for PLC2

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

3.2.6 Application

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

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

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

Property 1: PLC1_MemoryRegisters_MW0, PLC1_MemoryRegisters_MW1,
PLC2_MemoryRegisters_MW0, PLC2_MemoryRegisters_MW1,
PLC3_MemoryRegisters_MW0 and PLC3_MemoryRegisters_MW1

registers contain constant integer values (40, 80, 10, 20, 0, 10 respectively)³. The authors speculate that they may be (relative) hardcoded **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.

Property 2: PLC1_Coils_QX01, PLC1_Coils_QX02, PLC2_Coils_QX01, PLC2_Coils_QX02, PLC3_Coils_QX01 and PLC3_Coils_QX03 contain mutable binary (Boolean) values. The authors speculate that these registers can be associated with the **actuators** of the system.

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

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

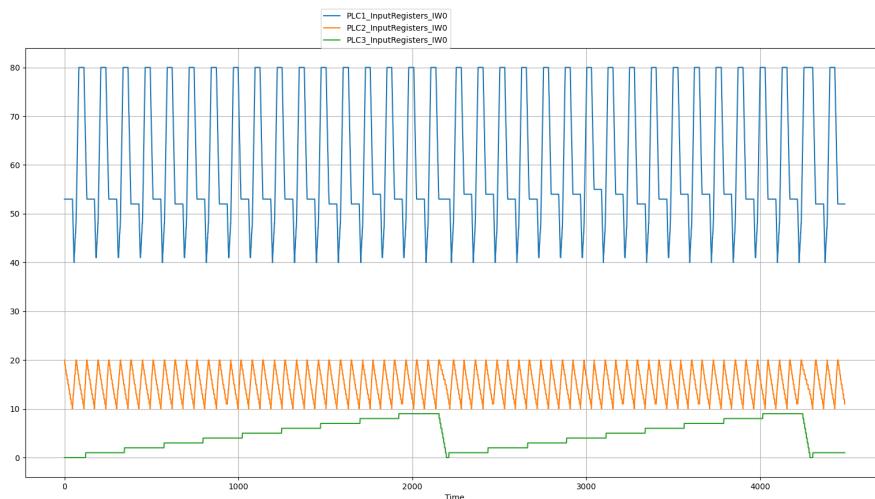


Figure 3.5: Execution traces of InputRegisters_IW0 on the three PLCs

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

Conjecture 2:

- the relative setpoints for PLC1_InputRegisters_IW0 are 40 and 80;

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

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

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

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

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

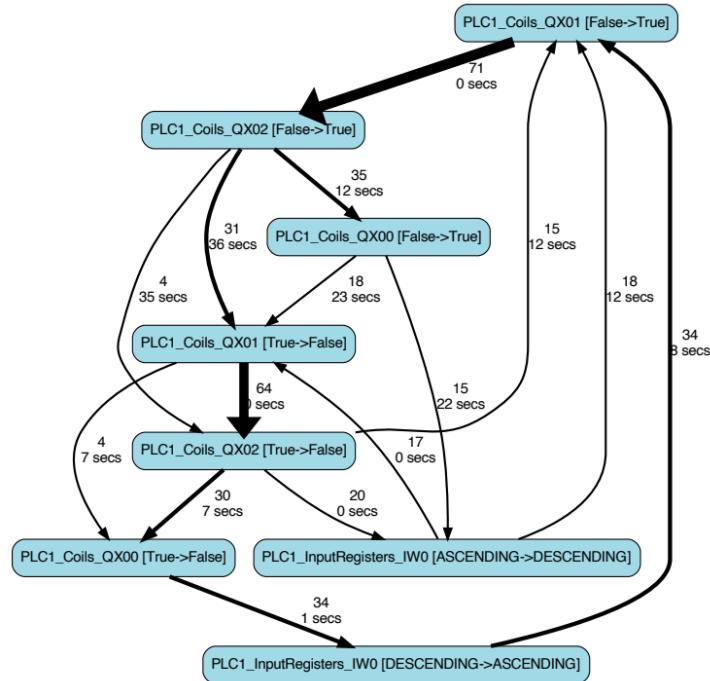
Property 4: PLC2 sends messages to PLC1 (see Figure 3.6b) which are recorded to PLC1_Coils_QX02.

Conjecture 3: PLC2_Coils_QX00 determines the trend in tank T-202 (Figure 3.6b).

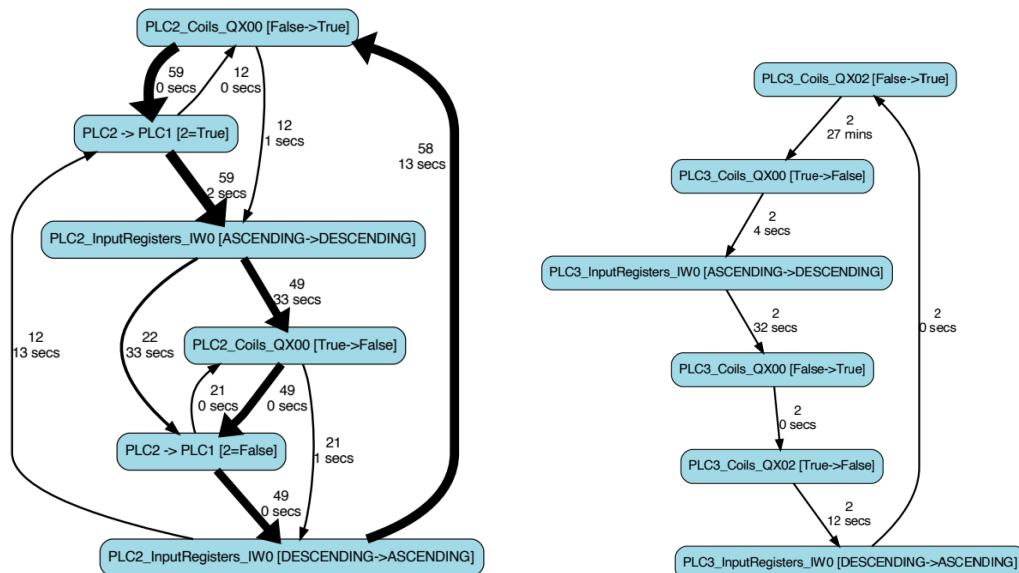
When this register is set to *True*, the input register PLC2_InputRegisters_IW0 related to the tank controlled by PLC2 starts an **ascending trend**; vice versa, when the coil register is set to *False*, the input register starts a **descending trend**.

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

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

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

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

However, it was possible to have confirmation of the conjectures made in the previous stages of the analysis: starting with the setpoints, analyzing the output of the invariants returned by Daikon⁴ reveals that

```
PLC1_InputRegisters_IW0 >= PLC1_MemoryRegisters_MW0 == 40.0  
PLC1_InputRegisters_IW0 <= PLC1_MemoryRegisters_MW1 == 80.0  
PLC2_InputRegisters_IW0 >= PLC2_MemoryRegisters_MW0 == 10.0  
PLC2_InputRegisters_IW0 <= PLC2_MemoryRegisters_MW1 == 20.0  
PLC3_InputRegisters_IW0 >= PLC3_MemoryRegisters_MW0 == 0.0  
PLC3_InputRegisters_IW0 <= PLC3_MemoryRegisters_MW1 == 9.0
```

i.e., that the `MemoryRegisters_MW0` and `MemoryRegisters_MW1` registers of each PLC contain the **absolute minimum and maximum setpoints**, respectively (*Property 5*).

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

`PLC1_Coils_QX01 == PLC1_Coils_QX02 == PLC2_Coils_QX00`

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

Regarding the **relationships between actuator state changes and measurement trends**, invariant analysis yields the results summarized in the following rules:

Property 7: Tank T-202 level *increases* iff `PLC1_Coils_QX01 == True`. Otherwise, if `PLC1_Coils_QX01 == False` will be *non-increasing*.

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

`PLC2_InputRegisters_IW0 == PLC2_MemoryRegisters_MW0 == 20.0 && PLC2_slope > 0`

is verified. On the opposite hand, if the coil is *False*, the condition

`PLC2_InputRegisters_IW0 == PLC2_MemoryRegisters_MW0 == 20.0 && PLC2_slope <= 0` is verified. The *slope* is increasing if > 0 , decreasing if < 0 , stable otherwise.

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

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

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

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

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

3.2.7 Limitations

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

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

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

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

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

does not help the methodology to be adaptable to different systems communicating with different protocols (sometimes even multiple ones on the same system).

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

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

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

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

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

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

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

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

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

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

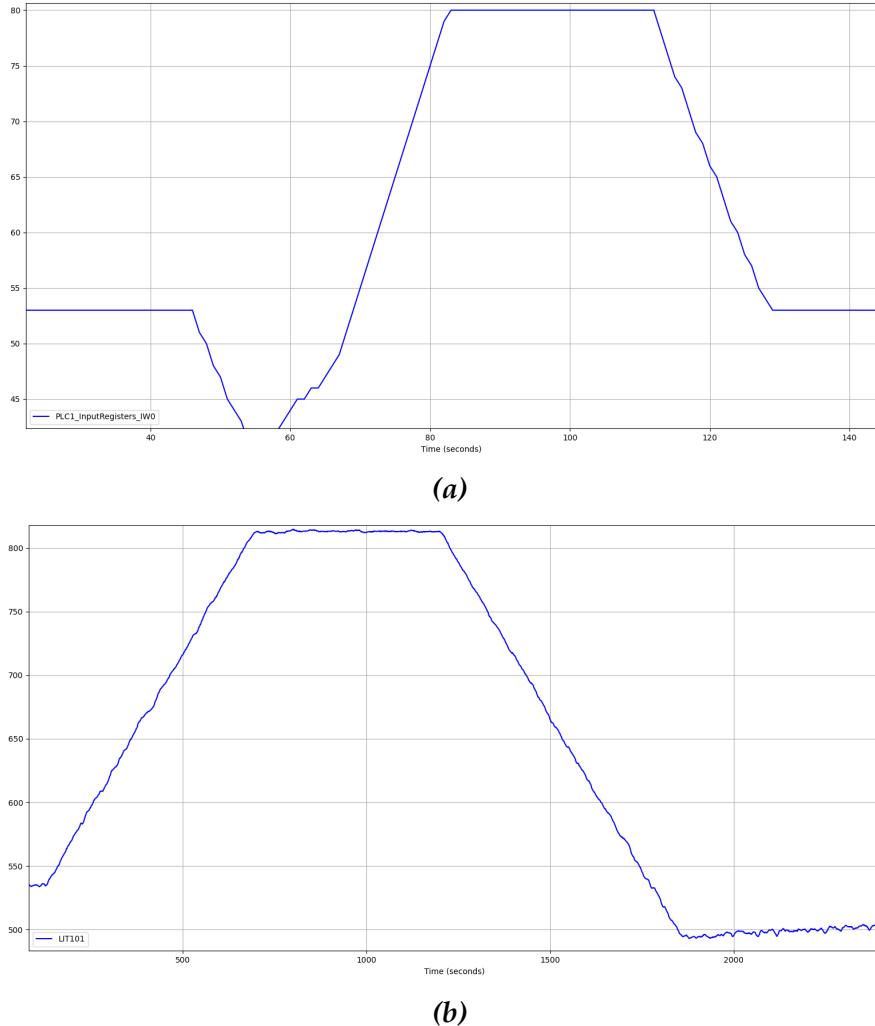


Figure 3.7: Water physics compared: simulated physics in the Simulink model (a) and physics in a real system (iTrust SWaT) (b). Fluctuations in the tank level in (b), almost completely absent in (a), can be appreciated.

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

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

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

First limitation Merging the datasets of all individual PLCs into a single dataset representing the entire system can be a sound practice if the system to be analyzed is (very) small as is the testbed analyzed here, consisting of a few PLCs and especially a few registers. If, however, the complexity of the system increases, this type of merging can become counterproductive and make it difficult to analyze and understand the data obtained in subsequent steps.

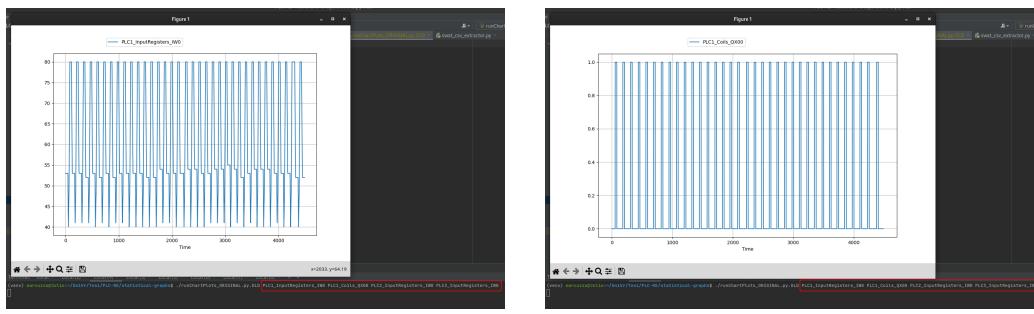
In short, there is no possibility to analyze only a subsystem and thus make the analysis faster and more understandable. Moreover, a data gathering can take up to days, and the analyst/attacker may need to make an analysis of the system isolating precise time ranges, ignoring everything that happens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

Second limitation Regarding the additional attributes, looking at the code of the script that performs the enrichment, we observed that **some attributes were manually inserted** after the merging phase: we are referring in particular to the attributes PLC x _Max_safety and PLC x _Min_safety, whose references were moreover hardcoded into the script, and the *slope* whose calculation method we mentioned in the previous paragraph about the testbed limitations.

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

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

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

First limitation While displaying charts for individual registers still provides useful information about the system such as the distinction between actuators and measurements and the general trend of the latter, single display does not allow one to catch, or at least makes it difficult, the relationship that exists between actuators and measurements, where it exists, because a view of the system as a whole is missing.

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

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

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

ical analysis (except maybe something marginal); the second observation pertains to the presentation of statistical information obtained from the histogram plot. In certain cases, the histogram plot itself can overshadow the displayed statistical information. These statistics are actually shown on the terminal from which the script is executed. However, to an inattentive or unfamiliar user, these statistics may be mistaken for debugging output or warnings, as they coincide with the display of the histogram plot window, which takes the focus (see Figure 3.9).

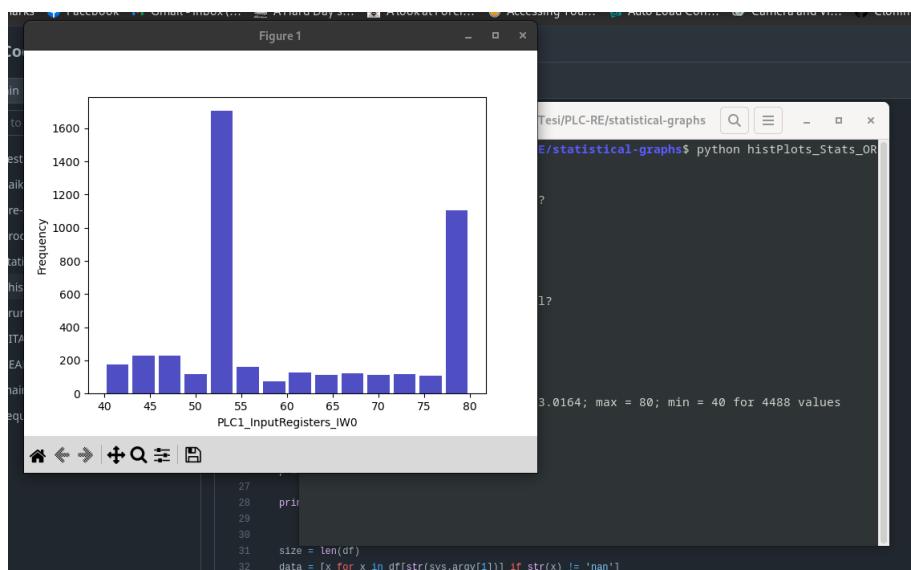


Figure 3.9: Histogram plot overshadowing statistical information shown on the terminal window in the background

In general, however, little statistical information is provided.

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

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

First limitation In a real system, network communications are much more numerous and involve many more devices even in the same second: finding the exact correspondence with what is happening in the cyber physical system becomes much more difficult.

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

Second limitation The other limiting aspect of the business process mining phase is the **process mining software** used to generate the activity diagram. As mentioned in Section 3.2.5, the process mining software used by Ceccato et al. is **Disco**: this is commercial software, with an academic license lasting only 30 days (although free of charge), released for Windows and MacOS operating systems only, which makes its use under Linux systems impossible except by using emulation environments such as Wine.

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

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

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 of invariants, one invariant per line, many of no use and without correlating invariants that may have common features or deriving

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

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

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

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

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

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

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

Extending and Generalizing Ceccato et al.'s Framework

IN Chapter 3, we presented the state of the art of *process comprehension* of an Industrial Control System (ICS) with a focus on the methodology proposed by Ceccato et al. [1][Section 3.2], explaining what it consists of, its practical application on a testbed, and most importantly highlighting its limitations and critical issues (see Section 3.2.7).

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

and transparent.

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

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

- **Section 4.1** will introduce a **novel framework** that we have developed to address the limitations of Ceccato et al.’s framework. Section will provide a brief overview of the framework’s features and structure;
- in **Section 4.2** the **features of our framework** will be demonstrated through their application to the different steps of the methodology. Practical examples will be used to illustrate the functionality of the framework. To facilitate this, we will utilize a more complex and detailed testbed compared to the one employed by Ceccato et al.

4.1 The Proposed New Framework

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

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

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

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

Single Programming Language The implementation of Ceccato et al.'s tool involved the use of different programming languages for each of the phases, ranging from Python to Java, and even including Bash scripting.

However, we believe that this heterogeneity introduces challenges in terms of user-friendliness and intuitiveness. It becomes more difficult for users to navigate and operate the tool effectively. Furthermore, the utilization of multiple technologies complicates code maintenance and the addition of new features, especially when managed by a single person who may have expertise in only one language while being less familiar with others.

Considering these factors, we have made the decision to adopt a single programming language for the framework to ensure homogeneity, simplify usability, and facilitate code maintenance for future users. Python has been chosen as the language of choice due to its simplicity, readability, versatility, and powerfulness. Additionally, Python benefits from a vast ecosystem of libraries and packages that cater to various requirements, making it a suitable choice for our framework.

System Independence One of the significant limitations we identified in Ceccato et al.'s tool, as highlighted in Section 3.2.7, is its **strong dependence on the specific testbed** being used. This means that the tool lacks the flexibility to configure parameters for analyzing different industrial systems.

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

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

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

To settle these issues, we enhanced the command-line interface in the proposed framework by adding new options and parameters. These new features provide the user with greater flexibility, enabling to specify parameters and options that allow for more in-depth analysis and focused results analyzing data more effectively and efficiently. With these enhancements, the framework has become more user-friendly, reducing the learning curve and making it more accessible to a wider range of users.

This, in turn, makes the framework more valuable and useful, increasing its adoption and effectiveness across a range of industrial control systems and applications.

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

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

4.1.1 Framework Structure

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

```
.
|-- config.ini
|-- daikon
|   |-- Daikon_Invariants
|   |-- daikonAnalysis.py
|   |-- runDaikon.py
|-- network-analysis
|   |-- data
|   |-- networkAnalysis.py
|   |-- export_pcap_data.py
|   |-- swat_csv_extractor.py
|-- pre-processing
|   |-- mergeDatasets.py
|   |-- system_info.py
|-- process-mining
|   |-- data
|   |-- process_mining.py
|-- statistical-graphs
|   |-- histPlots_Stats.py
|   |-- runChartSubPlots.py
```

Listing 4.1: Novel Framework structure and Python scripts

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

well as options related to specific analysis phases.

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

- [PATHS]: defines general paths such as the project root directory;
- [PREPROC]: includes parameters needed for the **pre-processing phase**, like the directory containing the raw datasets of the individual PLCs and *granularity* for the slope calculation. The granulariy is given in terms of the time interval over which the slope is calculated;
- [DATASET]: defines settings and parameters used during the **dataset enrichment** stage, for example the additional attributes;
- [DAIKON]: defines parameters needed for **invariant analysis** with Daikon, e.g. directories and files containing the outcomes of the analysis;
- [MINING]: contains parameters used during the **process mining** phase, such as data directory;
- [NETWORK]: includes specific settings for extracting the data obtained from the packet sniffing phase on the ICS network and converting it to CSV format. It also defines the **network protocols** that are to be analyzed.

4.1.2 Python Libraries and External Tools

As the framework has been developed entirely in Python, the objective was to minimize reliance on external tools and instead integrate various functionalities within the framework itself. The aim was to make the framework independent from external software. The only remaining external tool from the Ceccato et al. tool is Daikon. This choice was made

because there is currently no better alternative or Python package available that offers the same functionalities as Daikon.

Instead, the framework extensively utilizes Python libraries for handling various functionalities and input data. The core libraries on which the framework relies are:

- **Pandas**, also used in the Ceccato et al.'s tool for dataset management, but whose use here has been deepened and extended
- **NumPy**, often used together with Pandas to perform some operations to support it;
- **Matplotlib**, for managing and plotting graphical analysis;
- other scientific libraries such as **SciPy**, **StatsModel** [60] and **NetworkX** [61], for mathematical, statistical and analysis operations on the data;
- **GraphViz**, for the creation of activity diagrams in the process mining phase.

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

4.2 Analysis Phases

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

- **Section 4.2.1:** this section will present the **testbed** used to demonstrate examples of the framework's application;
- **Section 4.2.2:** the introduction of the new **network traffic analysis** phase (*Phase 0*) will be discussed. This phase aims to understand the structure of the industrial system's network and analyze its communication 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.

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 [15]. 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.

The testbed comprises several components, including:

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

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

Despite its moderate complexity, this testbed provides an ideal platform for presenting straightforward and concise examples of the framework's behavior. It enables us to effectively demonstrate the potential of the framework and facilitate a deeper understanding of its underlying methodology.

4.2.2 Phase 0: Network Analysis

The objective of the network analysis presented in this section is to offer users valuable information regarding the communication process within an industrial control system and a broader perspective on network communications, delving into previously unexplored aspects of the system. These additional dimensions provide a deeper understanding of the system's behavior and characteristics. This analysis aims to provide users with an overview of the communication between PLCs at level 1 (see Figure 2.1), as well as the communication between PLCs and devices at higher levels such as HMIs and Historian servers (see Section 2.1 for ICSs architecture). This also allows us to have a better understanding of the system architecture and network topology. The analysis focuses on industrial protocols used and the information exchanged.

By reconstructing the network communication structure using data obtained from the network traffic sniffing process, users can gain a comprehensive understanding of the behavior of the underlying industrial sys-

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

4.2.2.1 Extracting Data from PCAP Files

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

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

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

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

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

- **-s or --singledir:** operates differently from the previous option mentioned. This option enables the extraction of data from each individual PCAP file within the specified directory. The extracted data is then saved in separate CSV datasets, which are stored in the directory specified by the `split_dir` directive in the `config.ini` file. This functionality proves useful when dealing with exceptionally large PCAP files, where merging them together for export might consume significant time and resources. By utilizing this option, the extraction procedure becomes lighter and more manageable. The extracted data in separate CSV files can be utilized in the later stages of the Network Analysis process;
- **-t or --timerange:** this functionality enables users to specify a specific time period within the PCAP files from which they wish to extract relevant information.

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

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

4.2.2.2 Network Information

During this stage, the exported CSV data is processed to derive valuable information regarding network communications and the structure of the network itself. The objective is to identify and establish relationships between IP addresses present on the network, thereby determining the sources and destinations of communications. Furthermore, the analysis detects the protocols used for each communication and quantifies the various types of requests made.

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

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

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

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

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

- **-s or --srcaddr:** allows for specifying the source IP address for which you wish to display the incoming and outgoing communications. By providing the source IP address as an argument, the script will focus on showcasing the communications associated with that particular IP address;
- **-d or --dstaddr:** enables the user to specify the destination IP address for which you want to display the incoming and outgoing communications. By providing the destination IP address as an argument, the script will concentrate on presenting the communications associated with that specific IP address.

The parameters related to IP addresses, including source and destination, are optional. It is possible to specify either one of them individually. For instance, if the user specifies only the source IP address, the script will display the network nodes with which it communicates on the outgoing side, along with the corresponding generated traffic. Similarly, if only the destination IP address is specified, the script will showcase the network nodes communicating with it on the incoming side, along with the relevant traffic data.

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

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

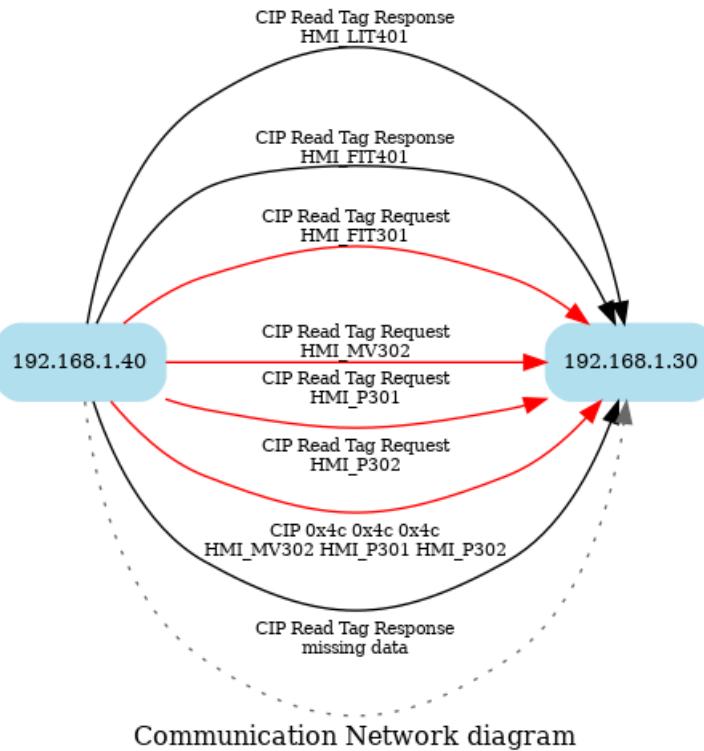


Figure 4.1: Network Communications between 192.168.1.40 (Source) and 192.168.1.30 (Destination)

The graph illustrates the network communications between the source IP address 192.168.1.40 and the destination IP address 192.168.1.30. Each arrow represents a communication between these IP addresses, showcasing the flow and direction of the interactions. The red arrows indicate requests initiated by the source IP address towards the destination IP, while the black arrows represent responses sent by the source IP in response to previous requests made by the destination IP. The gray dotted arrows represent responses for which the corresponding request is missing or unavailable for some reason. Overall, the graph distinguishes the different types of communications and provides insights into the request-response dynamics between the source and destination IP addresses. The graph is automatically generated and saved within the `$(project_dir)/network-analysis/data` directory.

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

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

4.2.3 Phase 1: Data Pre-processing

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

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

The combination of these limitations, along with the presence of hardcoded references to attributes and registers in the tool's code, makes the analysis of the system more challenging. Furthermore, it compromises the accuracy and reliability of the obtained results in terms of both quantity and quality.

In the proposed framework, these issues have been addressed by incorporating new features.

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

Secondly, we have **revamped the process of enriching** the resulting dataset by eliminating manual entry of additional attributes. Instead, users now have the flexibility to determine the type of additional attribute to associate with a specific register.

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

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

4.2.3.1 Subsystem Selection

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

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

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

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

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

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

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

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

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

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

4.2.3.2 Dataset Enrichment

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

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.

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

```
[DATASET]
timestamp_col = Timestamp
max_prefix = max_
min_prefix = min_
max_min_cols_list = lit|ait|dpit
prev_cols_prefix = prev_
prev_cols_list = mv[0-9]{3}|p[0-9]{3}
trend_cols_prefix = trend_
trend_cols_list = lit
trend_period = 150
slope_cols_prefix = slope_
slope_cols_list = lit
```

Listing 4.3: config.ini parameters for dataset enriching

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

timestamp_col denotes the name of the column in the dataset that holds 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.

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

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

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

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

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

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

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

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

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

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

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

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

or conditions.

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

However, if perturbations significantly affect the measurement readings, calculating the slope on individual time intervals may result in an inaccurate trend definition, irrespective of the chosen granularity. In such cases, the fluctuating nature of the measurements due to perturbations can introduce errors in the slope calculation, making it less reliable as an indicator of the actual trend.

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

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

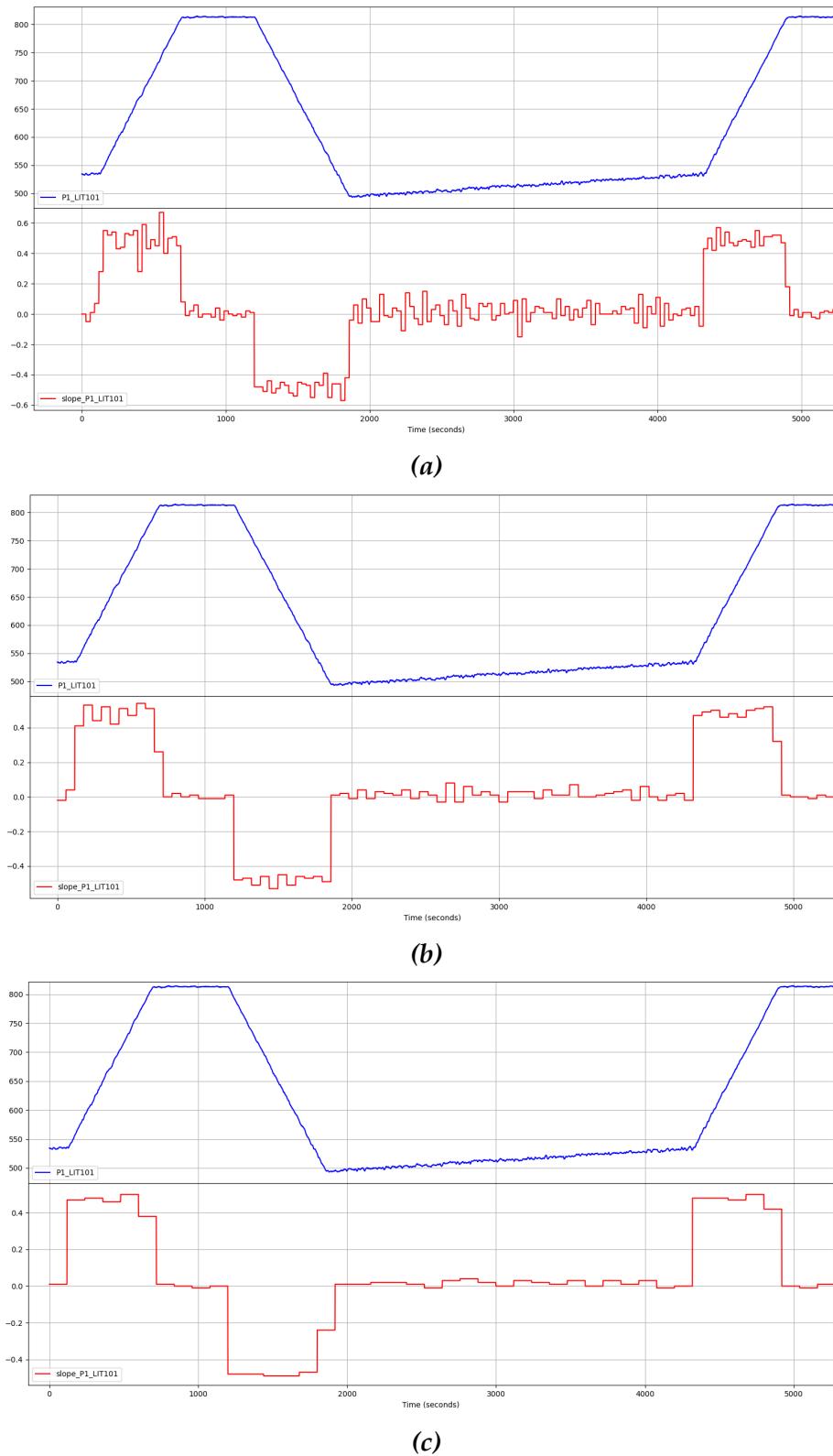


Figure 4.3: Slope comparison with granularity 30 (a), 60 (b) and 120 seconds (c)

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

The solution to this problem involves applying techniques to reduce the "noise" in the data, aiming to achieve a more linear trend in the measurement curve. By minimizing the effects of perturbations, we can calculate slopes more accurately.

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

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

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

Line simplification algorithms [64] aim to reduce the complexity of a polyline or curve by approximating it with a simplified version composed of fewer points. By selectively removing redundant or less significant points, line simplification algorithms help reduce storage space and computational requirements while preserving the overall shape and characteristics of the original line.

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

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

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

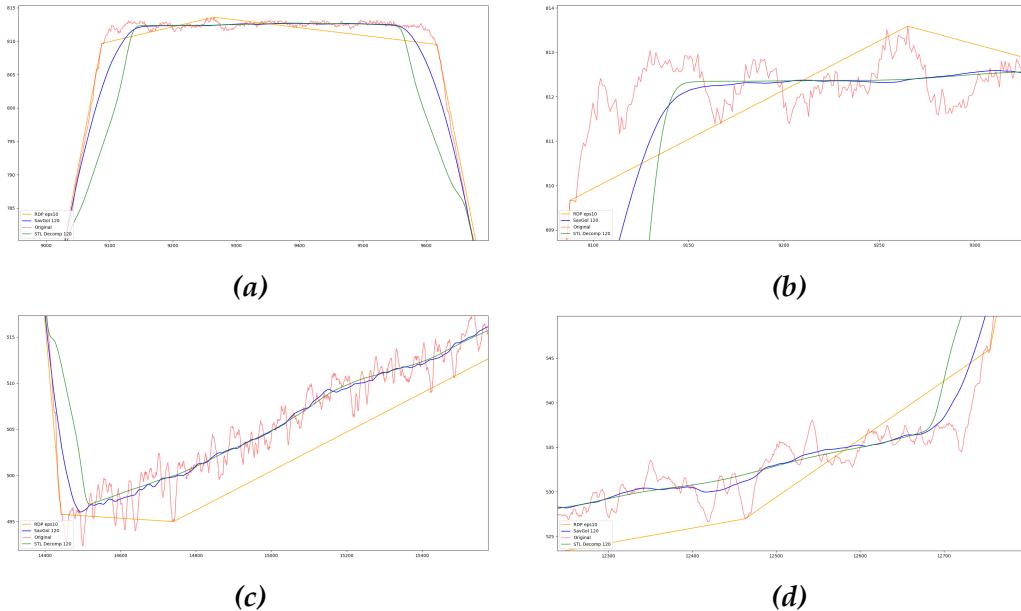


Figure 4.4: Savitzky-Golay filter (blue line), STL decomposition (green) and RDP algorithm (orange) comparison

By applying the STL decomposition, we observe a notable enhancement in slope calculation even when using a low granularity. Figure 4.5 demonstrates that, with the same granularity as shown in Figure 4.3a, the slope values, albeit exhibiting fluctuations, consistently align with the underlying trend of the data curve. The introduced lag resulting from the decomposition's periodicity is responsible for the observed delay.

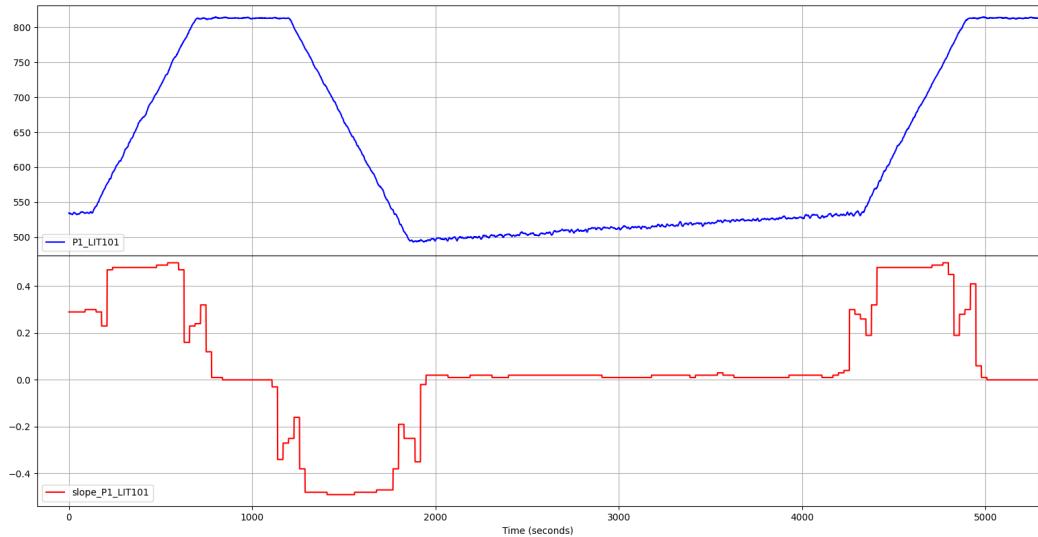


Figure 4.5: Slope after the application of the STL decomposition

The periodicity, which defines the sampling time window for decomposition and the level of noise smoothing, can be configured using the `trend_period` directive in the `config.ini` file.

During the slope calculation, the analysis will be performed on the data from the additional measurement trend attributes specified in the `trend_cols_list` directive of the configuration file, rather than on the original unfiltered data.

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

4.2.3.3 Datasets Merging

During this step, the datasets of the individual PLCs are merged, resulting in two separate datasets. The first dataset is enriched with additional attributes but excludes the timestamp column. This dataset is in-

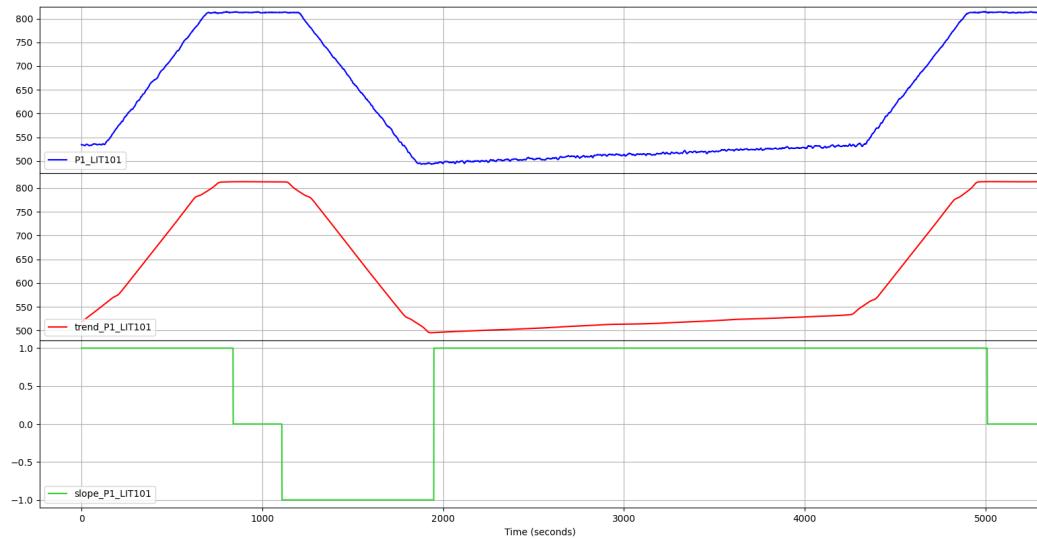


Figure 4.6: The new slope representation (green line) and the smoothed measurement data obtained with the STL decomposition (red)

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

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

4.2.3.4 Preliminary Analysis of the Obtained Subsystem

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

Relying on an analysis based on a combination of Daikon and Pandas this script performs a quick analysis of the dataset allowing to **estimate**, albeit approximately, the **type of registers** (sensors, actuators, ...), also identifying possible maximum and minimum values of measurements and hardcoded setpoints. Furthermore, leveraging the use of the additional attribute `prev_`, the `system_info.py` script is capable of deriving measurement values corresponding to state changes of individual actuators. This allows for the identification of specific measurements associated with the activation or deactivation of certain actuators within the system.

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

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

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

Actuators:
MV101 [0.0, 1.0, 2.0]
P101 [1.0, 2.0]

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

```
LIT101 {'max_lvl': 815.1, 'min_lvl': 489.6}
```

```
Hardcoded setpoints or spare actuators:  
P102 [1.0]
```

```
Actuator state changes:
```

LIT101	MV101	prev_MV101
800.7170	0	2
499.0203	0	1
800.5992	0	2
498.9026	0	1
800.7170	0	2
499.1381	0	1
801.3058	0	2
498.4315	0	1
801.4628	0	2
498.1567	0	1

LIT101	MV101	prev_MV101
805.0741	1	0
805.7414	1	0
805.7806	1	0
805.1133	1	0
804.4068	1	0

LIT101	MV101	prev_MV101
495.4483	2	0
497.9998	2	0
495.9586	2	0
495.8016	2	0
494.5847	2	0

LIT101	P101	prev_P101
536.0356	1	2
533.3272	1	2
542.1591	1	2
534.8581	1	2
540.5890	1	2

```

LIT101      P101      prev_P101
813.0031      2          1
813.0031      2          1
811.8256      2          1
812.7283      2          1
813.3171      2          1

Actuator state durations:
MV101 == 0.0
9   9   10   9   9   10   9   9   10   9

MV101 == 1.0
1174  1168  1182  1160  1172

MV101 == 2.0
669   3019  3012  3000  2981

P101 == 1.0
1073  1074  1061  1056  1056

P101 == 2.0
118   3133  3143  3125  3108

```

Listing 4.4: Example of preliminary system analysis

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

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

4.2.4 Phase 2: Graphs and Statistical Analysis

The introduction of the new *graph analysis* is motivated by from the requirement to provide users with a comprehensive overview of the (sub)system derived from the preceding pre-processing phase. The objective is to facilitate the identification of register types, enhance the understanding of relationships, and effectively grasp the dynamics among registers controlled by one or more PLCs. This analysis component serves to validate initial conjectures made during the preliminary analysis described in the previous section or generate new insights with the aid of visual chart representations. It enables users to gain a deeper understanding of the system by leveraging the support of visual charts.

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

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

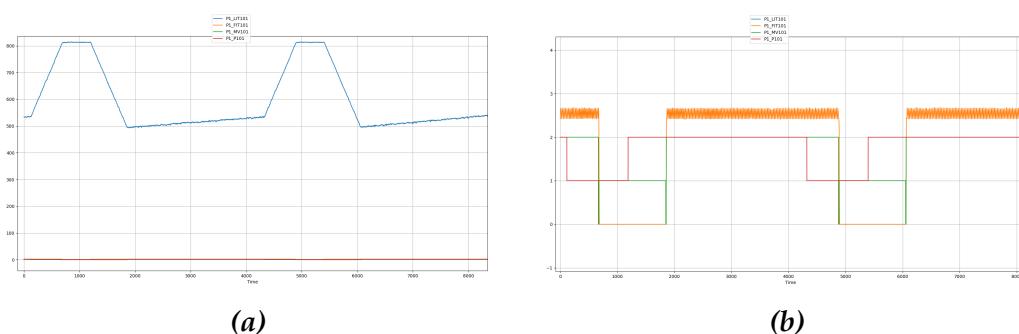
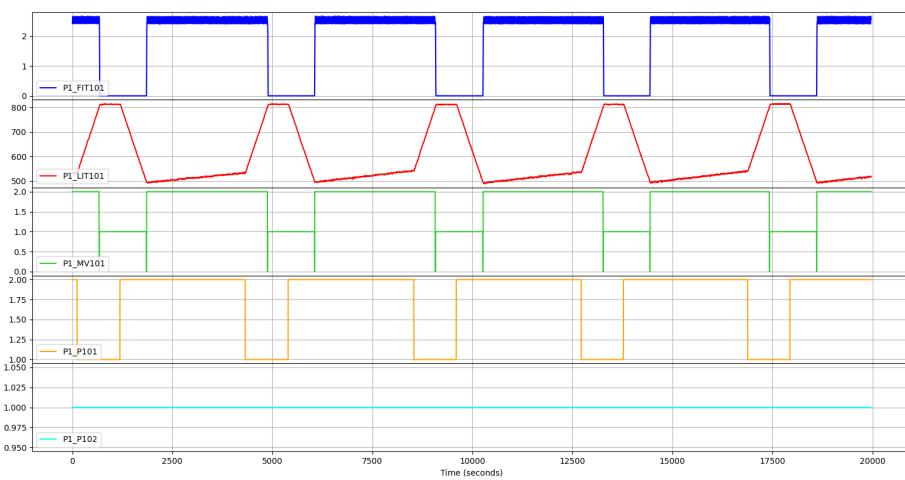


Figure 4.7: Plotting registers on the same y-axis

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.

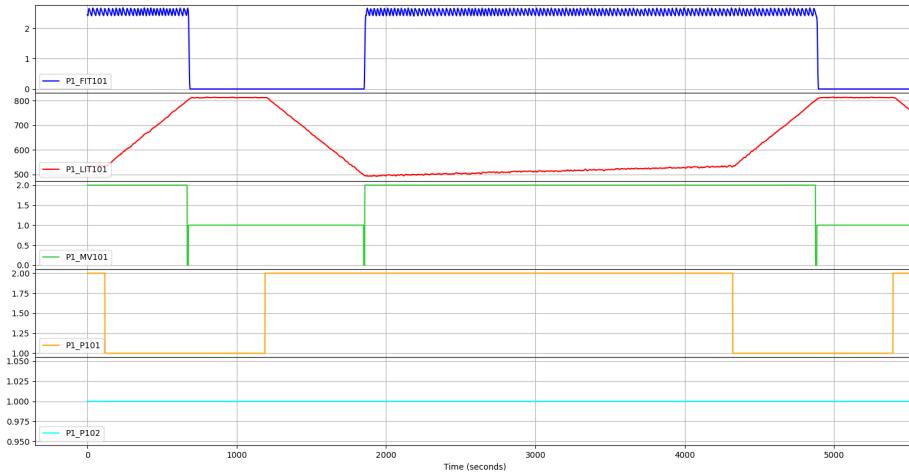
The solution to this issue is simple and effective: the use of **subplots**. Basically, each register corresponds to a subplot of the graph that shares the time axis (the x-axis) with the other subplots, but keeps the y-axis of the values of each register independent. This maintains the readability and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

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

The script accepts the following command-line parameters:

- **-f or --filename:** specifies the CVS format dataset to read data from. The dataset must be within the directory containing the enriched datasets for the invariant analysis phase. If subsequent parameters are not specified, the script will display all registers in the dataset, excluding any additional attributes;
- **-r or --registers:** specifies one or more specific registers to be displayed;
- **-a or --addregisters:** adds one or more registers to the default visualization. This option is useful in case an additional attribute such as slope is to be analyzed;

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

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

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

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

4.2.5 Phase 3: Invariant Inference and Analysis

The phase of invariant inference and analysis has undergone a redesign and improvement to offer the user a more comprehensive and easier approach to identify invariants. This has been achieved through the application of new criteria to analyze and reorganize the Daikon analysis results. The outcome of this is a more compact presentation of information that highlights the possible relationships among invariants.

The new design not only enables the identification of undiscovered aspects of the system behavior but also confirms the hypotheses made during the earlier stages of analysis. This step is now semi-automated, unlike before, and allows for the analysis of invariants on individual actuator states and their combinations.

4.2.5.1 Revised Daikon Output

To streamline the process of identifying invariants quickly and efficiently, it is necessary to revise the output generated by standard Daikon analysis. The goal is to create a more compact and readable format for the output.

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

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

In each section only a single invariant per row is shown, without relating it in any way to the others: this makes it difficult to identify significant invariants and any invariant chain that might provide much more information about the behavior of the system than the single invariant.

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

```
aprogram.point:::POINT
P102 == prev_P102
FIT101 >= 0.0
MV101 one of { 0.0, 1.0, 2.0 }
P101 one of { 1.0, 2.0 }
P102 == 1.0
max_LIT101 == 816.0
min_LIT101 == 489.0
slope_LIT101 one of { -1.0, 0.0, 1.0 }
[...]
LIT101 > MV101
```

```

LIT101 > P101
LIT101 > P102
LIT101 < max_LIT101
LIT101 > min_LIT101
[...]
MV101 < min_LIT101
MV101 < trend_LIT101
P101 >= P102
P101 < max_LIT101
[...]
=====
aprogram.point:::POINT; condition="MV101 == 2.0 && LIT101 <
    ↪ max_LIT101 - 16 && LIT101 > min_LIT101 + 15"
MV101 == prev_MV101
P102 == slope_LIT101
MV101 == 2.0
FIT101 > MV101
FIT101 > P101
FIT101 > P102
FIT101 > prev_P101
MV101 >= P101
MV101 >= prev_P101
P101 <= prev_P101
=====
aprogram.point:::POINT; condition="not(MV101 == 2.0 &&
    ↪ LIT101 < max_LIT101 - 16 && LIT101 > min_LIT101 + 15)
    ↪ "
P101 >= prev_P101
Exiting Daikon.

```

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

In the presented framework, the output is simplified to **two sections**: the *general section* and the section related to the *user-specified condition*. The section related to the negated condition is eliminated as it is not relevant in this context and could lead to potential misinterpretation. Moreover, the relationships between invariants will be emphasized by utilizing **transitive closures**: transitive closure of a relation R is another relation, typically 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 $R \subset R^+$ [69].

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

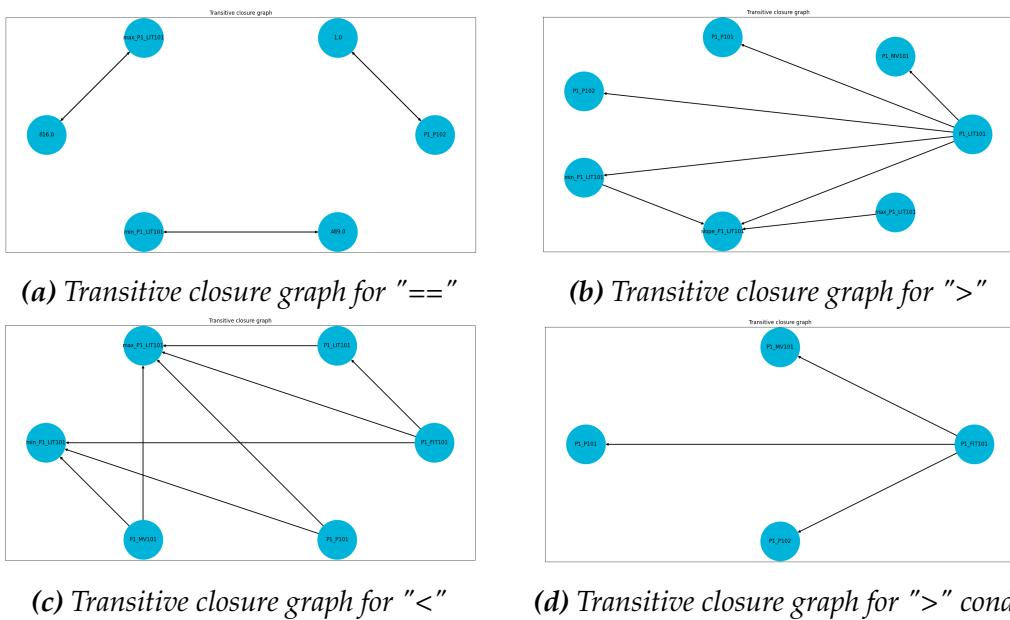


Figure 4.9: Example of transitive closure graphs for invariants in PLC1 of iTrust SWaT system

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

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

```

1 =====
2 General
3 =====
4 MV101 one of { 0.0, 1.0, 2.0 }
5 P101 one of { 1.0, 2.0 }
6 slope_LIT101 one of { -1.0, 0.0, 1.0 }
7 FIT101 != P101, P102
8 P102 == 1.0
9 max_LIT101 == 816.0
10 min_LIT101 == 489.0
11 LIT101 > MV101
12 LIT101 > P101
13 LIT101 > P102
14 LIT101 > min_LIT101 > slope_LIT101
15 FIT101 >= 0.0
16 P101 >= P102 >= slope_LIT101
17
18 =====
19 MV101 == 2.0 && LIT101 < max_LIT101 - 16 && LIT101 >
20   ↢ min_LIT101 + 15
21 =====
22 slope_LIT101 == P102
23 MV101 == 2.0
24 FIT101 > MV101
25 FIT101 > P101
26 FIT101 > P102
27 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.

4.2.5.2 Types of Analysis

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

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.

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.

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

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

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

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

Analysis on single actuator states Analysis on the states of individual actuators is the simplest: after the user is prompted to input the measurement, chosen from a list of likely available measurements, the script recognizes the likely actuators and the relative states of each, using the same mixed Daikon/Pandas technique adopted in the preliminary analysis during the pre-processing phase.

For each actuator and each state it assumes, a single Daikon analysis is performed, eventually placing the condition on the maximum and minimum level of the measurement.

The result of these analyses are saved in the form of text files in a directory having the name corresponding to the analyzed actuator and contained in the default parent directory

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

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

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

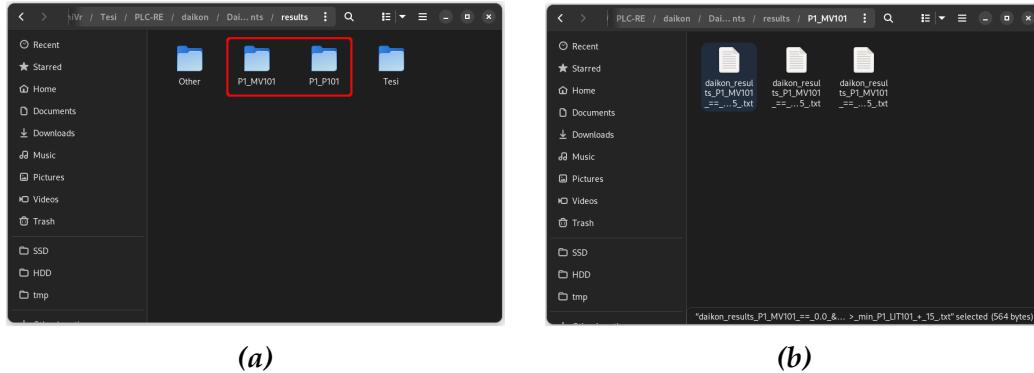


Figure 4.10: Directory (a) and outcome files (b) for the single actuator states analysis

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.

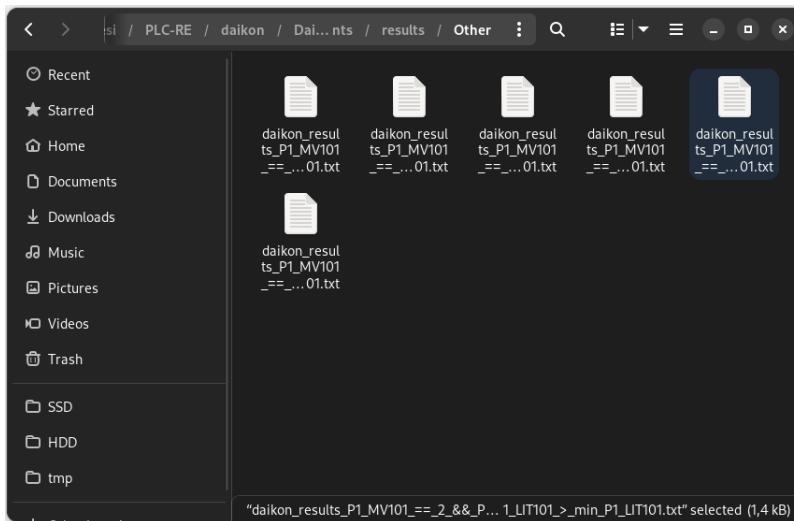


Figure 4.11: Daikon outcome files for system configuration analysis. Each file represents a single system state

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

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

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

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

```

1 =====
2 General
3 =====
4 MV101 <= P101 ==> FIT101 >= 0.0
5 MV101 <= P101 ==> MV101 one of { 0.0, 1.0, 2.0 }
6 MV101 <= P101 ==> P101 one of { 1.0, 2.0 }
7 MV101 <= P101 ==> slope_LIT101 one of { -1.0, 0.0, 1.0 }
8 MV101 > P101 ==> FIT101 > MV101
9 MV101 > P101 ==> FIT101 > P101
10 MV101 > P101 ==> FIT101 > P102
11 MV101 > P101 ==> FIT101 > slope_LIT101
12 MV101 > P101 ==> MV101 == 2.0
13 MV101 > P101 ==> MV101 > P102
14 MV101 > P101 ==> MV101 > slope_LIT101
15 MV101 > P101 ==> P101 == 1.0
16 MV101 > P101 ==> P101 == P102
17 MV101 > P101 ==> P101 == slope_LIT101
18 MV101 > P101 ==> slope_LIT101 == 1.0
19 [...]
20
21 =====
22 MV101 == 2 && P101 == 1 && LIT101 < max_LIT101 && LIT101 >
23   ↪ min_LIT101
24 =====
25 slope_LIT101 == P102 == P101 == 1.0
26 MV101 == 2.0
27 FIT101 > MV101

```

27 FIT101 > P101

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

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

Refining the Analysis In certain situations, the outcomes provided by the semi-automated analyses may not meet the user's expectations. For instance, the clarity of the slope value may be insufficient, or the user may wish to delve deeper into a specific aspect of the system to uncover additional invariants that were not previously identified. In such cases, the user has the option to conduct a more targeted and specific invariant analysis using the Python script `runDaikon.py`, which enables precise investigations of the system.

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

- **-f or --filename:** specifies the CSV format *enriched* dataset to read data from. Even in this case, the dataset must be located within the directory containing the enriched datasets;
- **-c or --condition:** specifies the condition for the analysis, which will be automatically overwritten in the `.sinfo` file. It is possible to specify more than one condition, but it is strongly recommended to use the logical operator `&&` to avoid undesired behaviors in Daikon outcomes;

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

During the execution of this analysis, a single output file is generated, containing the discovered invariants. The user can specify the directory where this file should be saved using the `-r` command-line option. By default, the file is stored in the directory

`$(project_dir)/daikon/Daikon_Invariants/results`. The output file serves as a record of the identified invariants and can be examined at a later time.

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

4.2.6 Phase 4: Business Process Analysis

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

In this updated Business Process, greater emphasis is placed on process mining related to the physical system. Our goal is to extract as much information as possible from the dataset, enabling us to promptly visualize the system's behavior and its various states. This approach allows us to validate the conjectures and hypotheses formulated in the previous phases, and potentially uncover hidden patterns that were previously undisclosed.

On the other hand, the aspect related to network communications was reconsidered to enable operation with multiple protocols, expanding beyond the limitations of Modbus.

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

4.2.6.1 Process Mining of the Physical Process

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

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

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

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

```
{  
    "MV101 == 2, P101 == 2": {  
        "start_value_LIT101": [  
            534.9366,  
            495.4483,  
            497.9998,  
            495.9586,  
            495.8016  
        ],  
        "end_value_LIT101": [  
            536.0356,  
            532.7384,  
            541.7273,  
            534.4656,  
            540.8245  
        ],  
        "slope_LIT101": [  
            0,  
            0.015,  
            0.018,  
            0.016,  
            0.018  
        ],  
        "trend_LIT101": [  
            "STBL",  
            "ASC",  
            "ASC",  
            "ASC",  
            "ASC"  
        ],  
        "time": [  
            119,  
            2464,  
            2477,  
            2452,  
            2440
```

```

    ],
    "next_state": [
        "MV101 == 2, P101 == 1",
        "MV101 == 2, P101 == 1"
    ]
},
"MV101 == 2, P101 == 1": {
    ...
}
"MV101 == 0, P101 == 1": {
    ...
}
...
}

```

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.

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

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

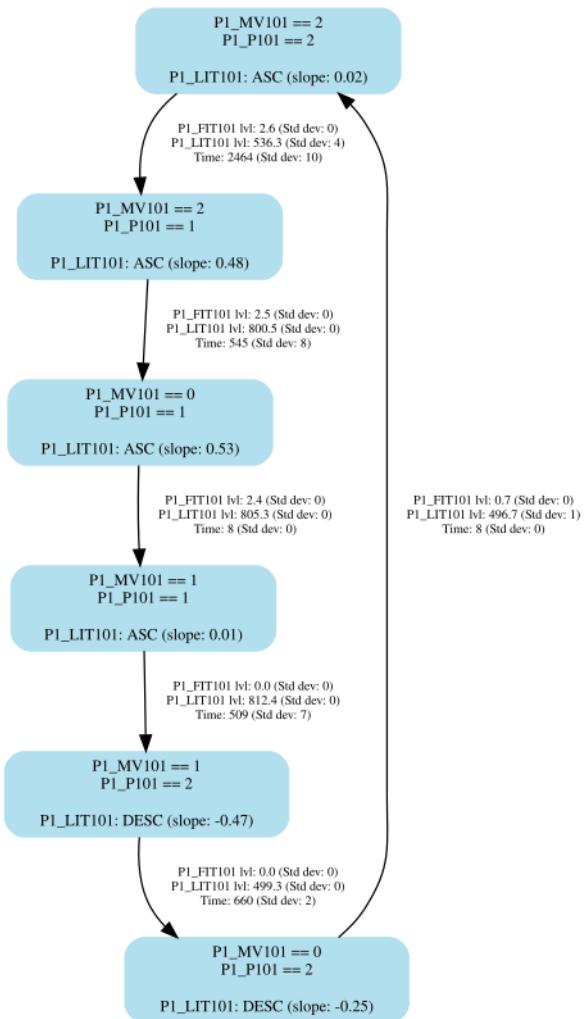


Figure 4.12: Activity diagram for PLC1 of the iTrust SWaT system

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

ity and the emphasis on the **temporal sequence of actuator states**. These insights might not have been clearly evident in earlier stages of the analysis. The diagram provides a visual representation that allows us to appreciate the recurring patterns and understand how the system evolves over time. By observing the transitions and relationships between different states, we gain a better understanding of the dynamics and behavior of the system.

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

4.2.6.2 Network Data

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

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

The analysis of the network data revealed that different protocols exhibit distinct behaviors when commanding changes in actuator states. Consequently, the previous approach proposed by Ceccato et al. becomes **impractical** when attempting to detect system state changes through network commands sent to the actuators.

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

4.2.7 Summary

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

Phase	Enhancements
General	- Independence from the analyzed system achieved through a general <i>config.ini</i> configuration file.
Improvements	- Command line functionality improved for enhanced usability. - Utilization of a single programming language throughout the implementation.

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

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

Table 4.1: Summary table of proposed framework enhancements

Case study: the iTrust SWaT System

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

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

5.1 Architecture

In contrast to the virtualized testbed discussed in Section 3.2.1 by Cecato et al., the iTrust SWaT system is composed entirely of physical hardware components. It encompasses various elements, starting from field

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

5.1.1 Physical Process

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

- P1. **taking in raw water:** feeds unfiltered water into the system
- P2. **chemical dosing:** adds chemicals to water useful for initial pretreatment;
- P3. **Ultra Filtration (UF) system:** the water is filtered through a semi-permeable membrane (ultrafiltration membrane) using the liquid pressure, effectively capturing impurities and suspended solids, as well as removing bacteria, viruses, and other pathogens present in the water;
- P4. **dechlorination:** removes residual chlorine from disinfected water using ultraviolet lamps;
- P5. **Reverse Osmosis (RO):** performs further filtration of the water;
- P6. **backwash process:** cleans the membranes in UF using the water produced by RO.

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

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

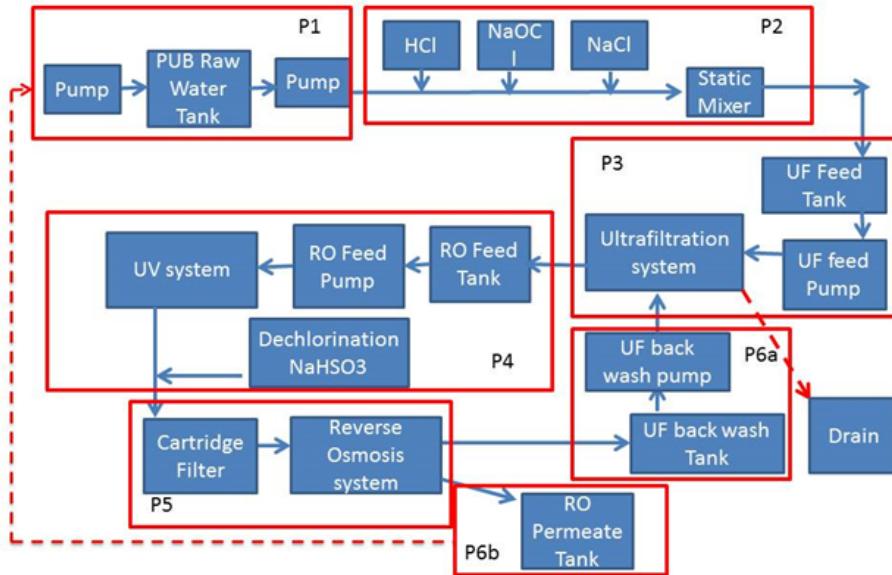


Figure 5.1: SWaT architecture

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

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

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

Figure 5.2: Sensors and actuators associated with each PLC

5.1.2 Control and Communication Network

The SWaT system's network architecture follows the principles of layering and zoning, which enable segmentation and control of traffic within the network.

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

- **Layer 3.5 – Demilitarized Zone (DMZ);**
- **Layer 3 – Operation Management (Historian);**
- **Layer 2 – Supervisory Control (Touch Panel, Engineering Workstation, HMI Control Clients);**
- **Layer 1 – Plant Control Network (PLCs) (Star Network);**
- **Layer 0 – Process (Actuator/Sensors and Input/output modules) (Ring Network).**

PLCs at Layer 1 communicate with their respective sensors and actuators at Layer 0 through a conventional ring network topology based on

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

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

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

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

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

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

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

 192.168.1.201

 Engineering Workstation

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

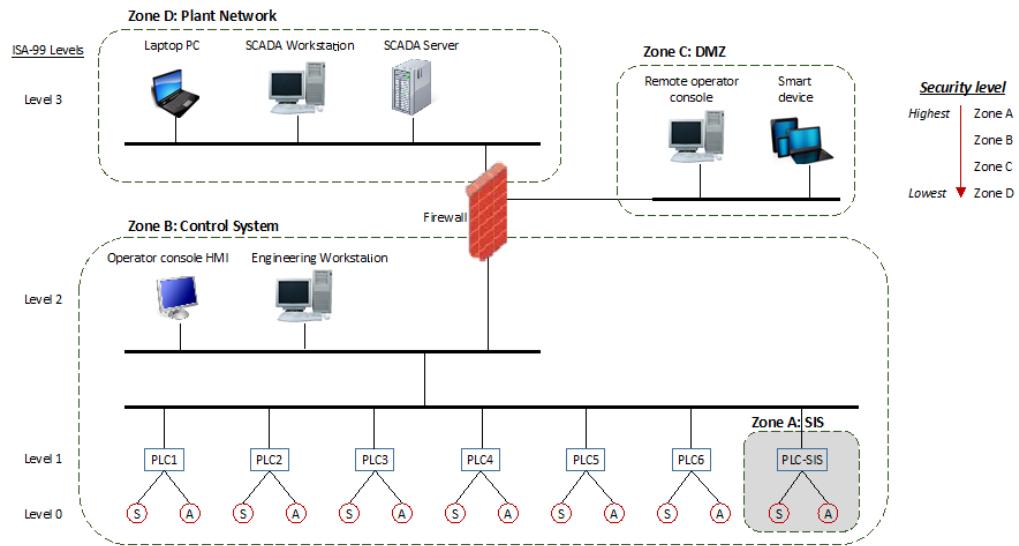


Figure 5.3: SWaT network architecture and zoning

5.2 Datasets

To facilitate the study and testing of technologies related to cyber security in Industrial Control Systems and critical infrastructure, iTrust offers researchers worldwide the opportunity to access a range of datasets [72]. These datasets consist of a collection of data obtained from the SWaT system, encompassing information on both the physical processes and network communications. The data is organized into different years, and researchers can request access to these datasets for their analysis and experimentation purposes.

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

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

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

5.2.1 Our Case Study: the 2015 Dataset

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

Physical Process Data The data collection process lasted 11 consecutive days, 24 hours per day. During the first 7 days, the system operated normally without any recorded attacks. However, attacks were observed during the remaining 4 days. The collected data reflects the impact of these attacks, leading to the creation of two separate CSV files: one containing the recorded data of SWaT during the system's regular operations, and the other containing data recorded during the days of the attacks. To ensure accurate information about the system, the dataset pertaining to the normal operations, which spans seven days, was chosen for analysis.

Data collection occurs at a frequency of one data point per second, 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).

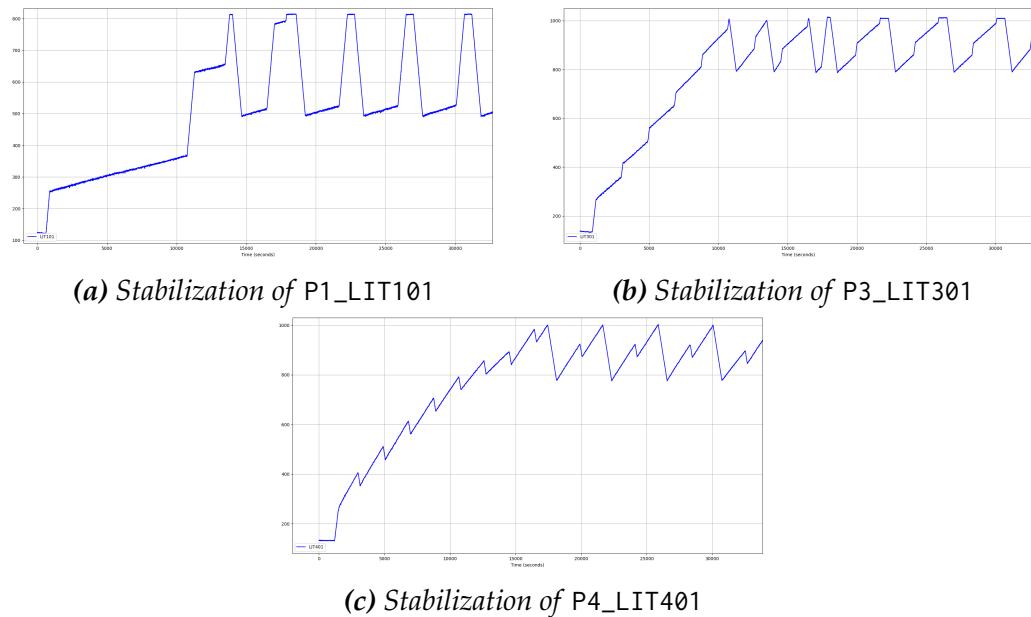


Figure 5.4: SWaT stabilization

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

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

To confirm that the communication protocol used is CIP over EtherNet/IP, we can examine the "Service / Destination Port" field. The destination port displayed is TCP port 44818, which corresponds to the TCP port commonly used for transporting **explicit messages** in the EtherNet/IP protocol (as explained in Section 2.2.6.2). Further evidence supporting the use of CIP protocol is observed in the *Modbus Function Code* field, consistently displaying a value of 76. In the Modbus protocol, there is no

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

Datasets for years beyond 2015 were not considered due to significant limitations that impede their usability. For instance, the 2017 dataset exhibits a high level of granularity in the physical system data, combined with short temporal periods for data gathering. As a result, valuable information is lost, and the network traffic data cannot be effectively utilized. On the other hand, post-2017 datasets suffer from the issue of being "spurious." This means that no distinction is made between the data collected during normal operations of the SWaT system and the data associated with system attacks, as was done in the 2015 dataset. Furthermore, the 2018 dataset lacks network traffic data altogether.

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

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

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

IN this chapter, our main objective is to apply the framework and methodology introduced in Chapter 4 to the case study of the iTrust SWaT system, as illustrated in Chapter 5. The purpose of this analysis is to assess the effectiveness and potential of the proposed framework within the context of a system that closely replicates a real-world water treatment plant, albeit on a smaller scale.

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

By focusing on selective components and leaving room for further exploration, we strike a balance between providing valuable insights and acknowledging the potential for additional research. This approach empowers the reader and interested individuals to explore the iTrust SWaT system further and leverage the proposed methodology and framework for a more comprehensive analysis.

6.1 Preliminary Operations

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

6.2 Planning the Analysis Strategy

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

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.

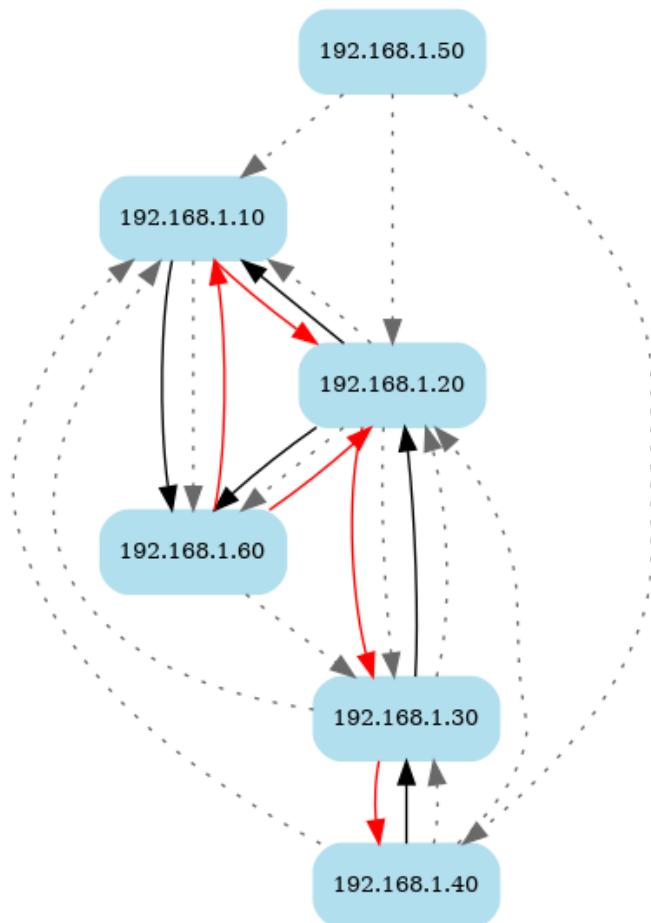


Figure 6.1: Simplified graph of the iTrust SWaT system network

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,

PLC6 communicates with both PLC1 and PLC2. On the other hand, the gray dotted arrows indicate communications for which we have knowledge of a response, but the corresponding request is unknown. For the purposes of our analysis strategy, we will not consider these communications within this context.

Based on our observations, the analysis strategy we will adopt involves considering sequential pairs of PLCs to effectively capture the relationships and implications between registers.

6.3 Reverse Engineering of the iTrust SWaT System

Before we delve into the analysis, it is important to provide some preliminary remarks.

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

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

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

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

About the Business Process Analysis In the end, the Business Process Analysis will focus solely on the physical process part. This is because the datasets of network traffic captures provided by iTrust for the year 2015 (as discussed in Section 5.2.1) are **incomplete**. While communications related to measurements are present, those associated with actuators are entirely missing, as well as additional communications related to other system characteristics that we observed in the datasets of subsequent years. As a result, we were unable to incorporate the network event recognition component into our Business Process Analysis. To implement this component, we would require complete and overlapping network data, along with a clean physical process dataset not affected by system attacks. Unfortunately, none of the available iTrust datasets fulfill these criteria.

6.3.1 Reverse Engineering of PLC1 and PLC2

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

6.3.1.1 Pre-processing - Preliminary Analysis

Measurements and Actuators Recognition Listing 6.1 shows the outcomes obtained from automatic recognition of likely measurements and actuators:

```

1  Actuators:
2  P1_MV101 [0.0, 1.0, 2.0]
3  P1_P101 [1.0, 2.0]
4  P2_MV201 [0.0, 1.0, 2.0]
5  P2_P203 [1.0, 2.0]
6  P2_P205 [1.0, 2.0]
7
8  Sensors:
9  P1_FIT101 {'max_lvl': 2.7, 'min_lvl': 0.0}
10 P1_LIT101 {'max_lvl': 815.1, 'min_lvl': 489.6}
11 P2_AIT201 {'max_lvl': 256.5, 'min_lvl': 252.9}
12 P2_AIT202 {'max_lvl': 8.4, 'min_lvl': 8.3}
13 P2_AIT203 {'max_lvl': 342.8, 'min_lvl': 320.0}
14 P2_FIT201 {'max_lvl': 2.5, 'min_lvl': 0.0}
15
16 Hardcoded setpoints or spare actuators:
17 P1_P102 [1.0]
18 P2_P201 [1.0]
19 P2_P202 [1.0]
20 P2_P204 [1.0]
21 P2_P206 [1.0]
```

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

Based on the results presented in Listing 6.1, the framework has identified P1_MV101, P1_P101, P2_MV201, P2_P203, and P2_P205 as **probable actuators**. The actuators denoted by the *Pxxx* notation are binary actuators

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

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

`P1_FIT101`, `P1_LIT101`, `P2_AIT201`, `P2_AIT202`, `P2_AIT203`, and `P2_FIT201` have been identified as **likely measurements**. Upon analyzing the range of values for register `P1_LIT101`, we observe a significant difference between the maximum and minimum values. This observation leads us to speculate that `P1_LIT101` could be identified as a **level sensor for the tank controlled by PLC1**. However, when examining registers `P1_FIT101`, `P2_FIT201`, `P2_AIT201`, and `P2_AIT202`, the small difference between their maximum and minimum values makes it unlikely that they represent additional tanks.

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

Some registers have been identified as **hardcoded setpoints or spare actuators** based on their constant values. These registers exhibit similarities to the previously recognized pump registers. It is plausible to speculate that these registers could correspond to **spare actuators**. Moreover, the constant value of 1 associated with these registers suggests that it may represent the **OFF state** of the pumps.

Actuator State Durations To gain a deeper understanding of the different states (0, 1, and 2) associated with valves `P1_MV101` and `P2_MV201`, we

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

```

1   Actuator state durations:
2   P1_MV101 == 0.0
3   9   9   10   9   9   10   9   9   10   9
4
5   P1_MV101 == 1.0
6   1174   1168   1182   1160   1172
7
8   P1_MV101 == 2.0
9   669   3019   3012   3000   2981
10
11  P2_MV201 == 0.0
12  8   8   8   9   9   8   9   9   9   9
13
14  P2_MV201 == 1.0
15  1057   1057   1045   1038   1039
16
17  P2_MV201 == 2.0
18  120   3135   3144   3127   3109

```

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

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

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

```

1   Actuator state changes:
2   ...

```

3	P1_LIT101	P1_P101	prev_P1_P101
4	536.0356	1	2
5	533.3272	1	2
6	542.1591	1	2
7	534.8581	1	2
8	540.5890	1	2
9
10	P1_LIT101	P1_P101	prev_P1_P101
11	813.0031	2	1
12	813.0031	2	1
13	811.8256	2	1
14	812.7283	2	1
15	813.3171	2	1
16

Listing 6.3: P1_P101 state changes in relation to P1_LIT101

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

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

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

Regarding the elements controlled by PLC2 and the sensor P1_FIT101,

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

6.3.1.2 Graphs and Statistical Analysis

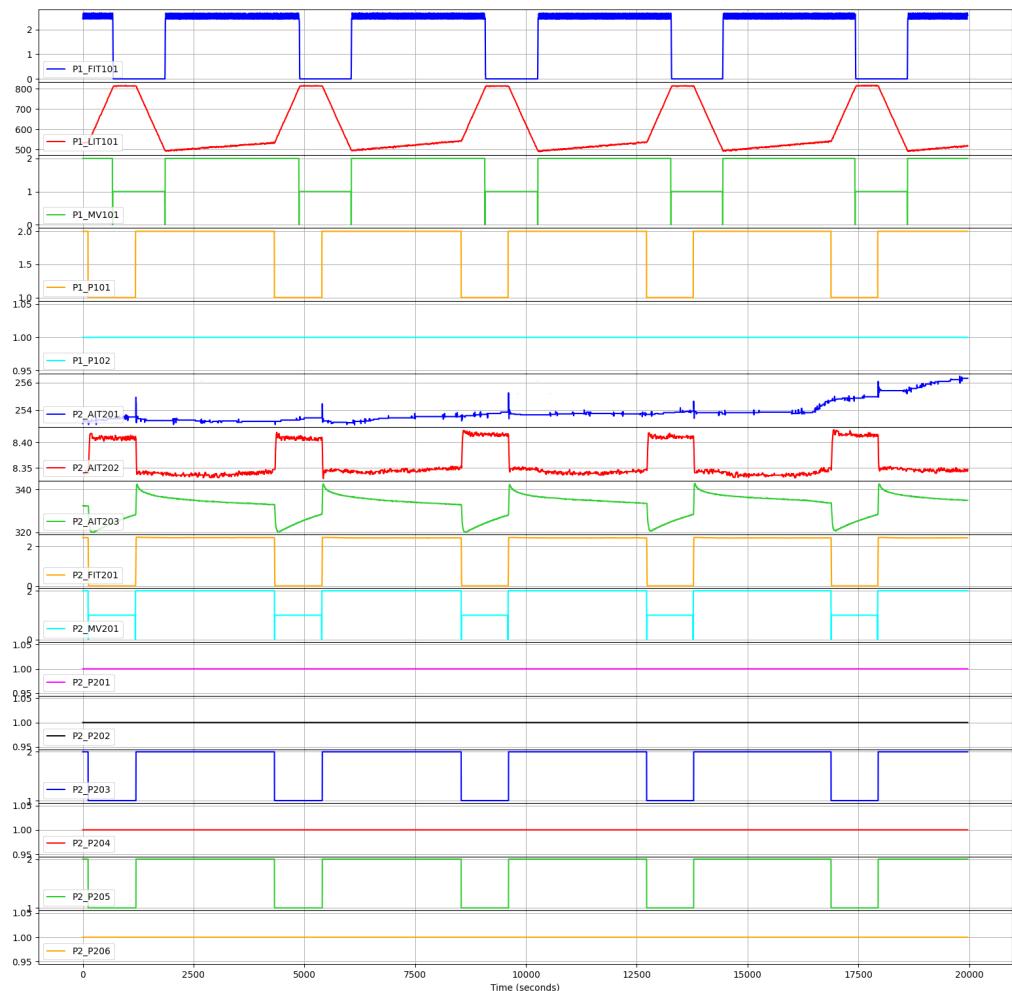


Figure 6.2: Chart of PLC1-2 registers

Figure 6.2 illustrates the graphical representation of the registers in

PLC1 and PLC2 and their respective trends.

The image provides additional support to the conjectures made during the preliminary analysis regarding the spare actuators. Furthermore, it is evident from the graphs that these spare actuators **do not appear to influence the trend** of any of the measurements. Therefore, based on this observation, we can confidently exclude these registers from further graphical analysis.

Figure 6.3 shows a clearer representation of the subsystem after removing the spare actuators.

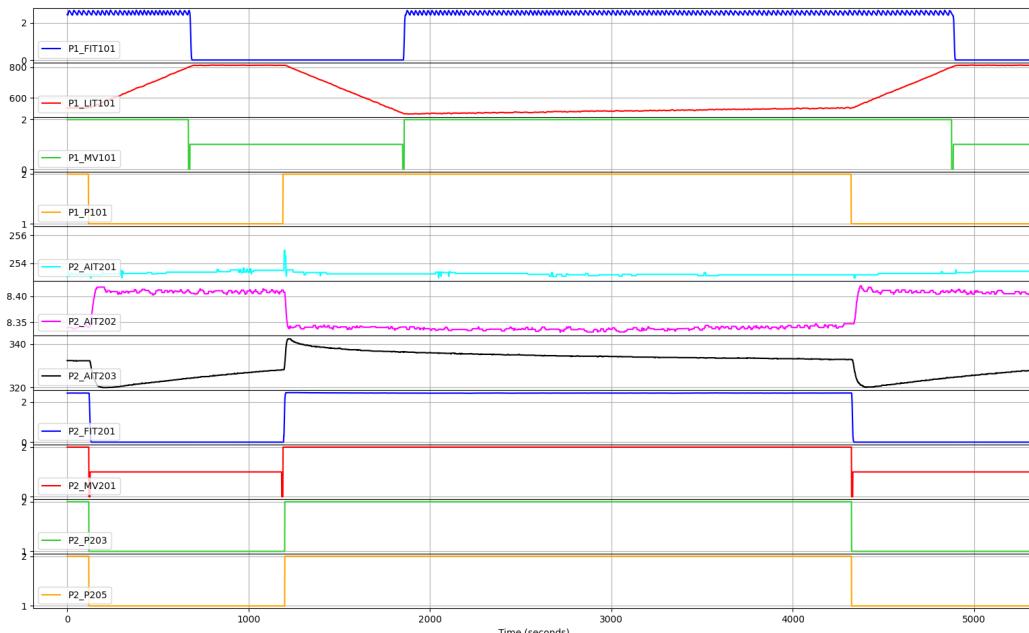


Figure 6.3: Chart of PLC1-2 registers without spare actuators (particular)

Figure provides furthermore additional insights that allow us to speculate on aspects that remained unexplained during the preliminary analysis. The most prominent aspect that stands out is the relationship between the level behavior of P1_LIT101 and the states of P1_P101 and P1_MV101. It appears evident that these two actuators **do not exhibit complementary behavior**, meaning that their states do not alternate in an ON-OFF pattern. 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:

```

1  P2_P206 == P2_P204 == P2_P202 == P2_P201 == P1_P102 == 1.0
2  P2_P205 == P2_P203
3  max_P1_LIT101 == 816.0
4  min_P1_LIT101 == 489.0
5  max_P2_AIT201 == 257.0
6  min_P2_AIT201 == 252.0
7  max_P2_AIT202 == 9.0
8  min_P2_AIT202 == 8.0
9  max_P2_AIT203 == 343.0
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.

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

```

1 =====
2 P1_MV101 == 1.0 && P1_LIT101 < max_P1_LIT101 - 16 &&
3   ↪ P1_LIT101 > min_P1_LIT101 + 15
4 =====
5 ...
6 P2_P205 == P2_P203 == P2_MV201 == P1_P101 == 2.0
7 P1_FIT101 == 0.0
8 slope_P1_LIT101 == -1.0
9 P2_FIT201 > P1_FIT101
10 P2_FIT201 > P1_MV101
11 P2_FIT201 > P1_P101
12 ...
13 =====
14 P1_MV101 == 2.0 && P1_LIT101 < max_P1_LIT101 - 16 &&
15   ↪ P1_LIT101 > min_P1_LIT101 + 15
16 =====
17 slope_P1_LIT101 == P1_P102
18 P1_FIT101 > P1_MV101
19 ...
20 P1_MV101 >= P1_P101
21 P1_MV101 >= P2_MV201
22 P1_MV101 >= P2_P203
23 P2_P203 >= P1_P101
24 ...

```

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

To prevent transient periods caused by water flow stabilization when

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

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

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

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

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

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

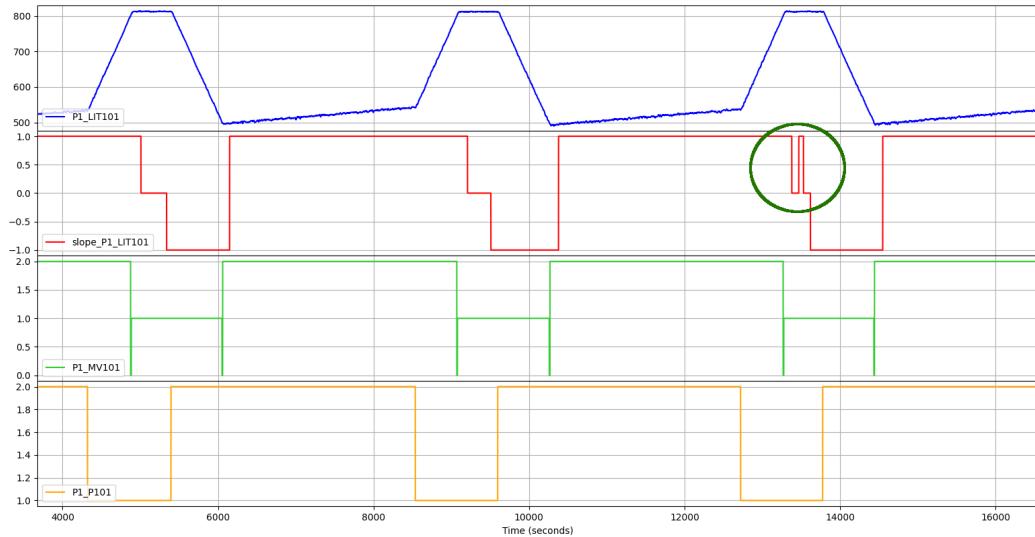


Figure 6.4: Slope calculation anomaly (in the circle)

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

6.3.1.4 Business Process Mining and Analysis

As explained in Section 4.2.6.1, the *process mining phase* applied to the physical system provides us with an immediate understanding of the system cycle and the chronological sequence of states. It enables us to determine the duration of time the system remains in a particular state and 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.

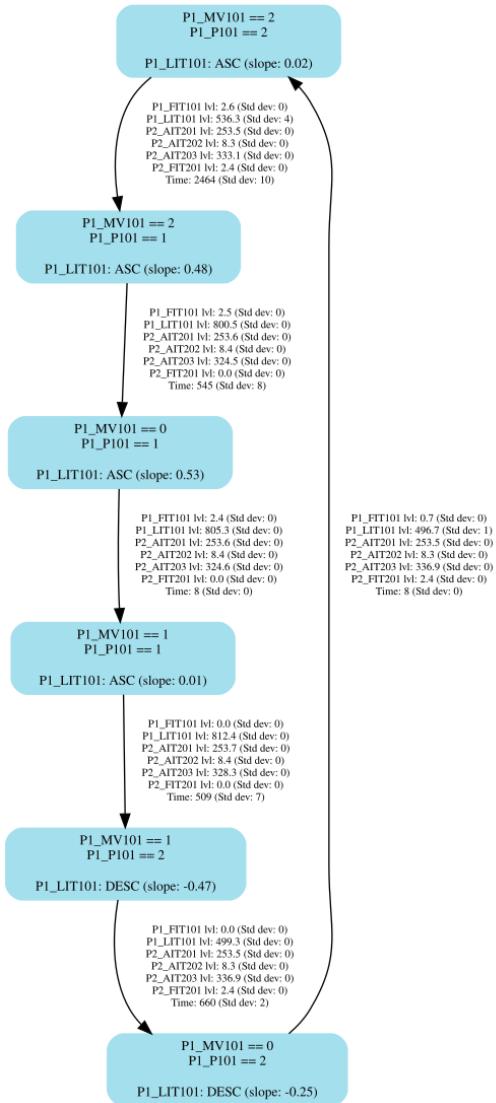


Figure 6.5: Activity diagram for PLC1-2

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

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

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

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

The data presented on the arcs in the activity diagram represents the measurement values relative to the time of system state changes, specifically the relative setpoints. These values are calculated based on the average data collected in each cycle. By analyzing this data, we can observe the tank level values at which the system undergoes configuration changes and how the trend of the tank level changes accordingly. Specifically, we can see that the trend changes from ascending to stable at a tank level

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

Furthermore, the data provides additional support for the hypothesis that the measurements associated with registers P2_AIT20x are not influenced by the tank's trends.

6.3.1.5 Properties

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

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

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

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

Table 6.1: Properties of the PLC1-2 subsystem

6.3.2 Reverse Engineering of PLC2 and PLC3

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

6.3.2.1 Pre-processing - Preliminary Analysis

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

```

1  Actuators:
2  ...
3  P3_MV301 [0.0, 1.0, 2.0]
4  P3_MV302 [0.0, 1.0, 2.0]
5  P3_MV303 [0.0, 1.0, 2.0]
6  P3_MV304 [0.0, 1.0, 2.0]
7  P3_P302 [1.0, 2.0]

8
9  Sensors:
10 ...
11 P3_DPIT301 {'max_lvl': 20.4, 'min_lvl': 0.0}
12 P3_FIT301 {'max_lvl': 2.4, 'min_lvl': 0.0}
13 P3_LIT301 {'max_lvl': 1014.5, 'min_lvl': 786.5}

14
15 Hardcoded setpoints or spare actuators:
16 ...

```

17 P3_P301 [1.0]

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

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

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

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

Actuator State Durations Let's proceed with the analysis of the actuator states' duration, as displayed in Listing 6.7. In this analysis, our focus will not be on examining the correspondence between values and the actual actuator states, as in Section 6.3.1.1. Instead, we will explore whether these actuators exhibit any distinct patterns or behaviors based on their duration.

```

1 Actuator state durations:
2 ...
3 P3_MV301 == 1.0
4 2527 4154 4154 4154 4094
5
6 P3_MV301 == 2.0
7 36 35 36 35 34
8
9 P3_MV302 == 1.0
10 662 138 654 138 656 139 658 137 656 137

```

```

11
12 P3_MV302 == 2.0
13 62 1783 1596 1787 1591 1791 1576 1803 1540 1782
14
15 P3_MV303 == 1.0
16 2526 4089 4088 4089 4028
17
18 P3_MV303 == 2.0
19 97 96 97 96 96
20
21 P3_MV304 == 1.0
22 689 1832 2206 1838 2203 1840 2191 1852 2152 1831
23
24 P3_P302 == 1.0
25 637 115 632 115 632 114 634 114 632 115
26
27 P3_P302 == 2.0
28 60 1821 1632 1825 1629 1829 1615 1841 1578 1820
29
30
31

```

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 observe that P3_P302, the pump in PLC3, exhibits a behavior

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

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

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

1	Actuator state changes:		
2	...		
3	P2_MV201	prev_P2_MV201	P3_LIT301
4	0	2	1000.2240
5	0	1	799.1140
6	0	2	1001.5060
7	0	1	799.1942
8	0	2	1001.5460
9	0	1	799.1140
10

Listing 6.8: P2_MV201 state changes in relation to P3_LIT301

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

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

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

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

6.3.2.2 Graphs and Statistical Analysis

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

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

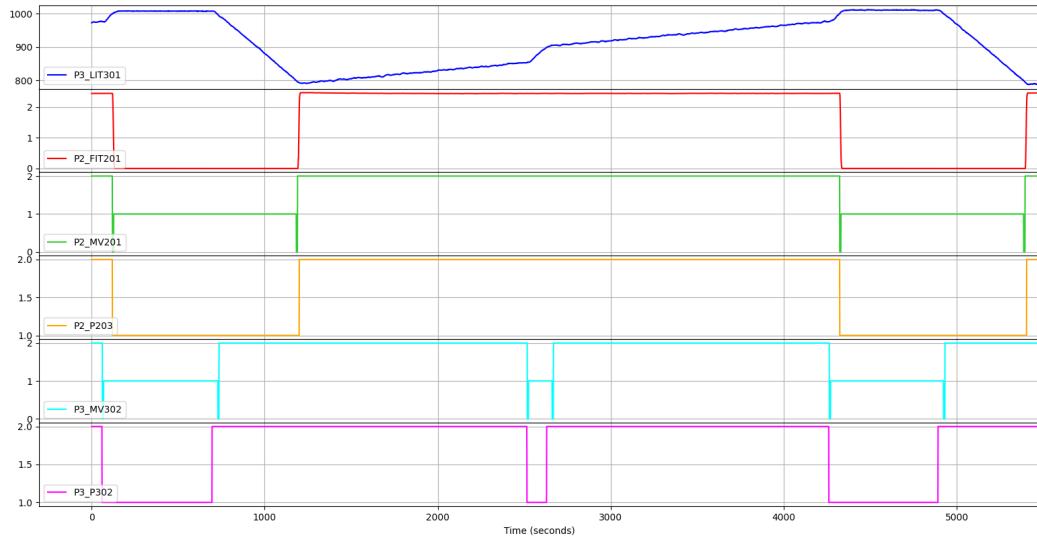


Figure 6.6: Verifying the conjecture about valve P2_MV201

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

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

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

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

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

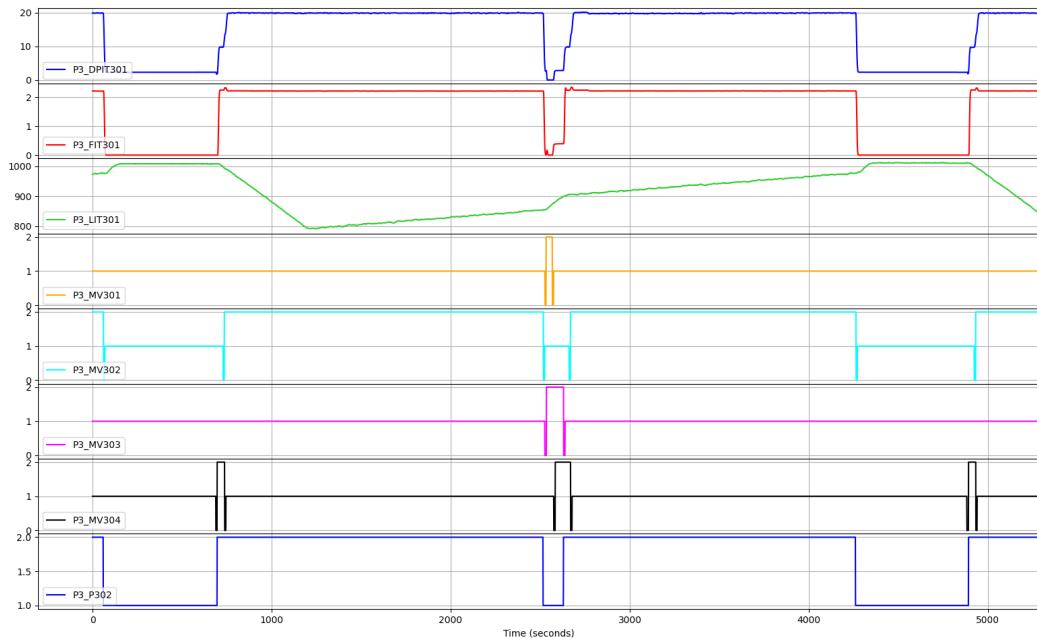


Figure 6.7: PLC3 registers

the incoming flow to another element outside the analyzed subsystem.

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

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

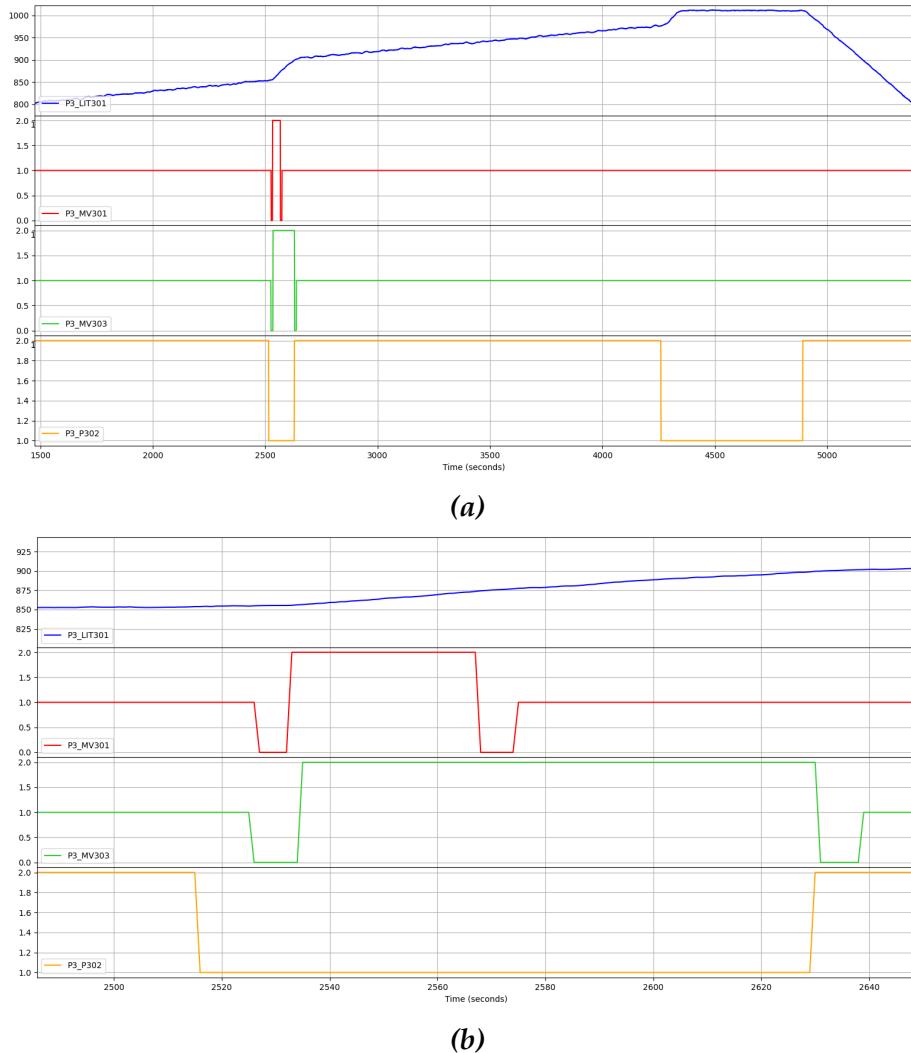


Figure 6.8: P3_MV301 and P3_MV303 analysis

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

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

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

6.3.2.3 Invariant Inference and Analysis

General Invariants The analysis of general invariants does not yield any significant insights. However, it confirms the maximum and minimum values of the measurements seen in Section 6.3.2.1 and identifies the presence of the spare actuator P3_P301.

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

```

1 =====
2 P2_MV201 == 1.0 && P3_LIT301 < max_P3_LIT301 - 20 &&
3   ↪ P3_LIT301 > min_P3_LIT301 + 24
4 =====
5 P3_P301 == P3_MV303 == P3_MV301 == P2_P205 == P2_P203 ==
6   ↪ P2_MV201 == 1.0
7 P3_P302 == 2.0
8 slope_P3_LIT301 == -1.0
9 P3_FIT301 > P2_MV201
10 P3_DPIT301 > P3_FIT301 > P3_P302
11 ...

```

```

10 =====
11 P2_MV201 == 2.0 && P3_LIT301 < max_P3_LIT301 - 20 &&
12   ↪ P3_LIT301 > min_P3_LIT301 + 24
13 =====
14 P2_P205 == P2_P203 == P2_MV201 == 2.0
15 slope_P3_LIT301 == P2_P201
16 P2_FIT201 > P2_MV201
17 P2_FIT201 > P2_P201
18 P2_FIT201 > P3_FIT301
19 ...
20 =====
21 P3_P302 == 1.0 && P3_LIT301 < max_P3_LIT301 - 20 &&
22   ↪ P3_LIT301 > min_P3_LIT301 + 24
23 =====
24 P2_P205 == P2_P203 == P2_MV201 == 2.0
25 slope_P3_LIT301 == P3_P302 == P2_P201
26 P2_FIT201 > P2_MV201
27 P2_FIT201 > P3_FIT301
28 ...
29 =====
30 P2_MV201 one of { 1.0, 2.0 }
31 slope_P3_LIT301 one of { -1.0, 1.0 }
32 P3_DPIT301 > P3_P302 > slope_P3_LIT301
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.

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

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

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

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

Unfortunately, no further useful information can be derived from this phase of the analysis. To address the remaining questions and clarify any outstanding issues, we will proceed to the next and final phase of the subsystem analysis.

6.3.2.4 Business Process Mining and Analysis

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

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

```

1   "slope_P3_LIT301": [
2     0.383, # from 970 to 1000
3     0.395, # from 850 to 900
4     0.384,
5     0.395,
6     0.354,
7     0.38,
8     0.388,
9     0.381,
10    0.385,
11    0.386
12  ],

```

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

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

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

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

- When the system is in the states [P2_MV201 == 2, P3_P302 == 2] and [P2_MV201 == 2, P3_P302 == 1], the level of the tank represented by P3_LIT301 is increasing. The tank level exhibits a faster growth rate in the latter state.
- When the system is in the state [P2_MV201 == 1, P3_P302 == 1], the tank level remains stable.
- When the system is in the state [P2_MV201 == 1, P3_P302 == 2], the tank level is decreasing.

The **absolute setpoints** for the tank level in subsystem PLC2-3 are defined as follows: the minimum setpoint is 800, the maximum setpoint is 1000, and there are additional setpoints at 850, 900, and 970, which correspond to specific changes in the slope of the tank level.

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

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

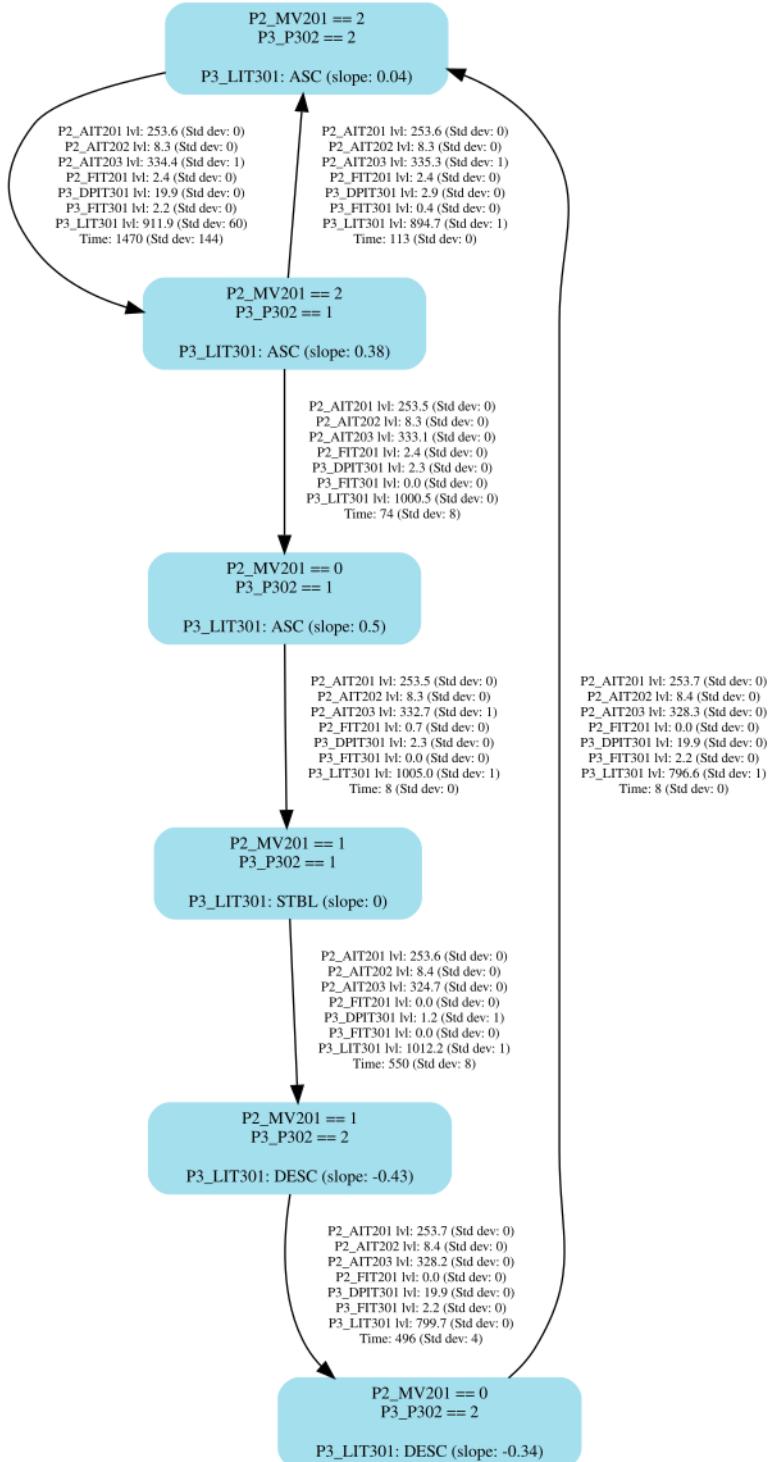


Figure 6.9: Activity diagram for PLC2-3

6.3.2.5 Properties

Table 6.2 provides a summary of the properties inferred from the conjectures made throughout the different stages of the analysis.

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

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

Table 6.2: Properties of the PLC2-3 subsystem

6.3.3 Reverse Engineering of PLC3 and PLC4

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

6.3.3.1 Pre-processing - Preliminary Analysis

Measurements and Actuators Recognition Listing 6.11 shows the outcomes obtained from automatic recognition of likely measurements and actuators for PLC4. We omit those related to PLC3 as they are already known.

```

1  Actuators:
2  ...

```

```
3
4 Sensors:
5 ...
6 P4_AIT401 {'max_lvl': 148.8, 'min_lvl': 148.8}
7 P4_AIT402 {'max_lvl': 191.1, 'min_lvl': 185.5}
8 P4_FIT401 {'max_lvl': 1.7, 'min_lvl': 1.7}
9 P4_LIT401 {'max_lvl': 1002.8, 'min_lvl': 775.8}
10
11 Hardcoded setpoints or spare actuators:
12 ...
13 P4_P401 [1.0]
14 P4_P402 [2.0]
15 P4_P403 [1.0]
16 P4_P404 [1.0]
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 P4_P404 are spare actuators, specifically pumps.

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

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

1	Actuator state changes:		
2	...		
3	P3_MV302	prev_P3_MV302	P4_LIT401
4	0	2	1000.5510
5	0	1	784.4911
6	0	2	922.1866
7	0	1	881.0818
8	0	2	1000.2820
9	0	1	786.1061
10	0	2	922.8018
11	0	1	881.8508
12

Listing 6.12: P3_MV302 state changes in relation to P4_LIT401

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

6.3.3.2 Graphs and Statistical Analysis

We will attempt to gain a deeper understanding of the pattern exhibited by the PLC4 registers by referencing Figure 6.10.

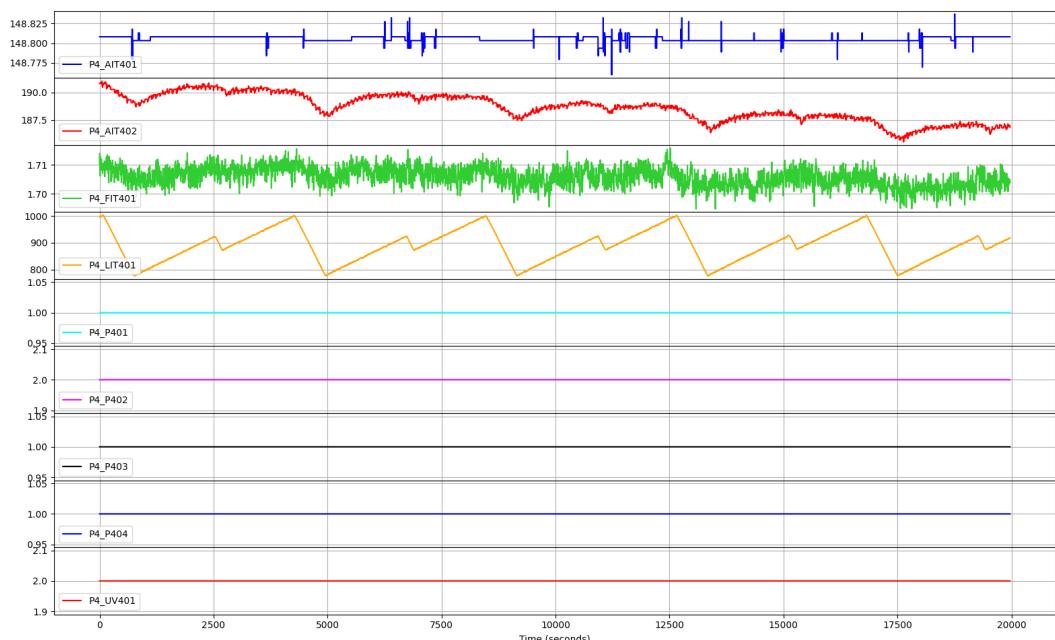


Figure 6.10: PLC4 registers

The image reveals interesting behavior in the P4_AIT401 and P4_AIT402 registers. Notably, P4_AIT401 exhibits a linear trend rather than a cyclic one, with values oscillating within a narrow range. This suggests that, similar to P2_AIT201 (refer to Section 6.3.2.2), this register may correspond to a sensor measuring a specific water property. On the other hand, P4_AIT402 appears to follow the level trend of sensor P4_LIT401, but with a downward cyclic pattern where each cycle starts at a lower level than the pre-

vious one. Given the limited value range and the similarity in naming conventions, it is highly likely that this register represents another water property sensor rather than a tank.

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

Now, let us refer to Figure 6.11 to seek confirmation regarding the conjecture that implicates valve P3_MV302 as the responsible actuator for tank filling.

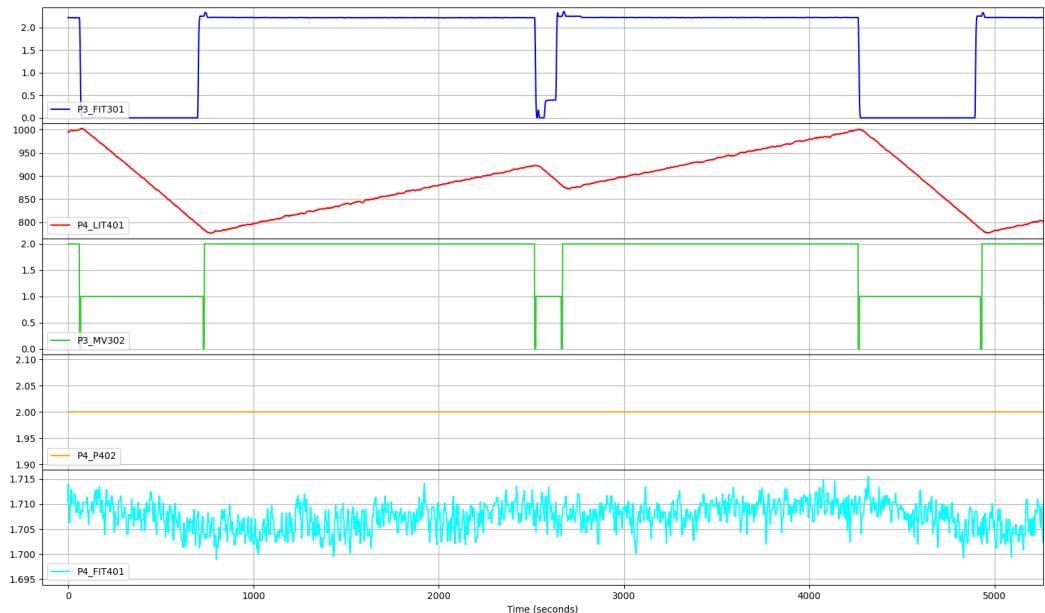


Figure 6.11: P3_MV302 and P4_LIT401 behaviors

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

Upon closer observation, it becomes apparent that the tank controlled by PLC4 **does not have plateau periods**. When the incoming water flow ceases, the tank immediately begins to empty. Based on our findings in

previous subsystems, we speculate that the actuator responsible for tank emptying could be the pump indicated by register P4_P402. This speculation is further supported by the nearly constant trend observed in sensor P4_FIT401.

Figure 6.12 depicts the correlation between the tanks within this subsystem and the actuators that are responsible for their filling cycle.

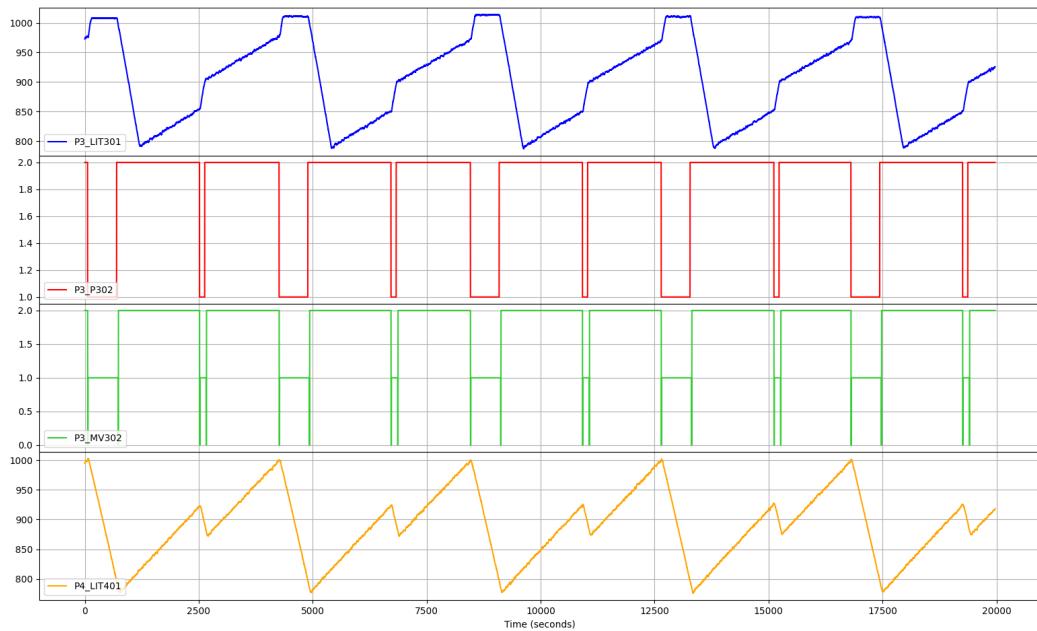


Figure 6.12: Tanks in subsystem PLC3-4 and their correlation.

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

6.3.3.3 Invariant Inference and Analysis

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

General Invariants Again, the analysis of the general invariants offers no new information compared to what we conjectured earlier. We therefore continue with the analysis of the current system configuration.

Analysis of the Current System Configuration The analysis of the current system configuration provides confirmation that when the system is in the state [$P3_MV302 == 1$, $P4_P402 == 2$], the tank level, as indicated by sensor $P4_LIT401$, shows a decreasing slope. Additionally, it is noted that the spare actuators align with the expected behavior.

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

```

1 =====
2 P3_MV302 == 2 && P4_P402 == 2 && P4_LIT401 < 990 &&
3   ↪ P4_LIT401 > 930
4 =====
5 ...
6 slope_P4_LIT401 == slope_P3_LIT301 == P4_P404 == P4_P403
7   ↪ == P4_P401 == P3_P301 == P3_MV304 == P3_MV303 ==
8   ↪ P3_MV301 == 1.0
9 ...

```

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

6.3.3.4 Business Process Mining and Analysis

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

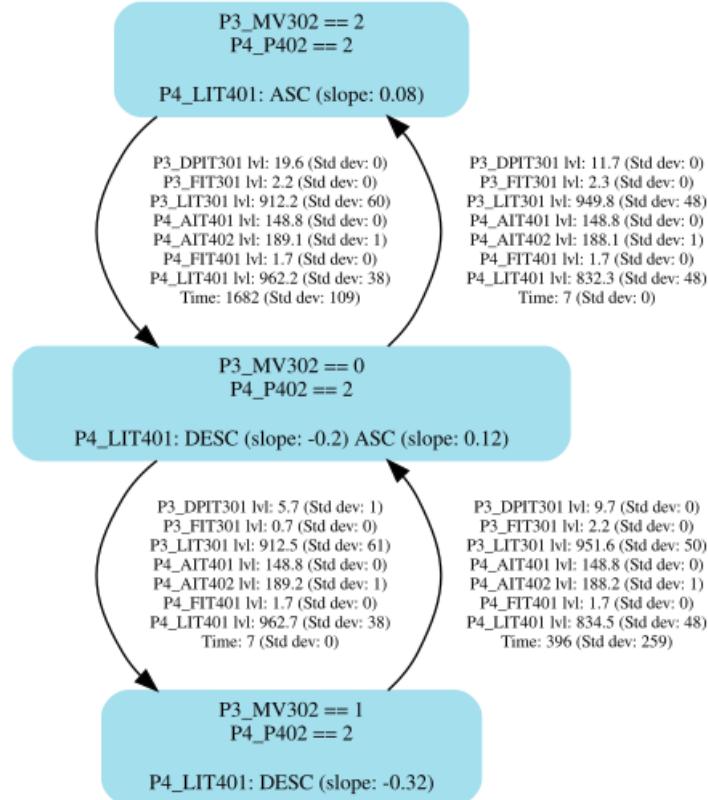


Figure 6.13: Activity diagram for PLC3-4

The diagram provides confirmation that when the system is in the state [P3_MV302 == 2, P4_P402 == 2], the tank level trend is increasing. Conversely, when the system is in the state [P3_MV302 == 2, P4_P402 == 2], the trend is decreasing.

Regarding the setpoints, based on the standard deviation of the data from P4_LIT401, the absolute setpoints for the system are as follows: 785 (minimum), 880 (relative minimum), 925 (relative maximum), and 1000 (maximum).

Furthermore, from the process mining phase, we can extract information about P3_FIT301 and P3_DPIT301. P3_FIT301 registers values greater than 2 when both pump P4_P402 and valve P3_MV302 are in the ON state, while it drops to an average of 0.7 when the valve is off. Similarly, P3_DPIT301

registers values close to 20 when the valve is ON, and these values decrease when the valve is OFF.

6.3.3.5 Properties

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

#	Statement	Derived from
P37	The registers P4_AIT401, P4_AIT402, P4_FIT401, and P4_LIT401 of PLC4 hold likely sensor measurements.	Preliminary Analysis Graphical Analysis
P38	The register P4_P402 of PLC4 holds likely actuator command. Its status is constantly ON.	Preliminary Analysis Graphical Analysis Business Process
P39	The registers P4_P401, P4_P403, and P4_P404 of PLC4 are associated to spare actuators.	Preliminary Analysis Graphical Analysis Invariant Analysis
P40	The register P4_LIT401 of PLC4 is associated to the level sensor of the tank controlled by PLC4.	Preliminary Analysis Graphical Analysis
P41	P3_MV302 and P4_P402 are the actuators responsible for the level behaviour of the water contained in the tank controlled by PLC4.	Graphical Analysis Invariant Analysis Business Process
P42	The level of the tank identified by the register P4_LIT401 increases when both P3_MV302 and P4_P402 are in the ON state. It decreases when P3_MV302 is in the OFF state.	Graphical Analysis Invariant Analysis Business Process
P43	The trend of the tank level controlled by PLC4 transition from ascending to descending when the level reaches approximately 925 and 1000. It changes from descending when the level reaches approximately 785 and 880.	Business Process
P44	Absolute setpoints are 785, 880, 925 and 1000	Business Process

P45	P4_FIT401 serves as a flow or pressure sensor to the P4_LIT401 register. It is related to the P4_P402 pump.	Graphical Analysis Business Process
P46	P4_P401 does not exhibit a cyclic trend.	Graphical Analysis
P47	The register P4_P402 exhibits a cyclic decreasing trend, which is closely linked to the trend observed in the P4_LIT401 register.	Graphical Analysis
P48	Both P4_AIT401 and P4_AIT402 serve as sensors for measuring certain properties of the water.	Preliminary Analysis Graphical Analysis
P49	The registers P3_MV301, P3_MV303, and P3_MV304 do not have an impact on the water level dynamics on the tank controlled by PLC4.	Graphical Analysis Business Process

Table 6.3: Properties of the PLC3-4 subsystem

6.4 PLCs Architecture

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

PLC	Physical device / Variable	Range	PLC register
PLC1	Level sensor for Tank 1	[489-815]	P1_LIT101
	Flow or pressure sensor	[0-2.7]	P1_FIT101
	Valve for Tank 1 inflow water	{0,1,2}	P1_MV101
	Pump for Tank 1 outflow water	{0,1}	P1_P101

	Spare actuator (OFF)	1	P1_P102
PLC2	Flow or pressure sensor	[0-2.4]	P2_FIT201
	Sensor for some water property	[252-256]	P2_AIT201
	Sensor for some water property	[8.3-8.4]	P2_AIT202
	Unkown	[320-342]	P2_AIT203
	Valve for Tank 3 inflow water	{0,1,2}	P2_MV201
	Spare actuator (OFF)	1	P2_P201
	Spare actuator (OFF)	1	P2_P202
PLC3	Pump (unknown role)	{0.1}	P2_P203
	Spare actuator (OFF)	1	P2_P204
	Pump (unknown role)	{0.1}	P2_P205
	Spare actuator (OFF)	1	P2_P206
	Level sensor for Tank 3	[768-1014]	P3_LIT301
	Flow or pressure sensor	[0-2.4]	P3_FIT301
	Flow or pressure sensor	[0-20.4]	P3_DPIT301
PLC3	Valve (unkonw role)	{0,1,2}	P3_MV301
	Valve for Tank 4 inflow water	{0,1,2}	P3_MV302
	Valve (unkonw role)	{0,1,2}	P3_MV303
	Valve (unkonw role)	{0,1,2}	P3_MV304
	Spare actuator (OFF)	1	P3_P301
	Pump for Tank 3 outflow water	[0,1]	P3_P302
	Level sensor for Tank 4	[775-1002]	P4_LIT401
PLC4	Flow or pressure sensor	[1.7]	P4_FIT401
	Sensor for some water property	[148.8]	P4_AIT401
	Sensor for some water property	[185-191]	P4_AIT402
	Spare actuator	1	P4_P401
	Pump for Tank 4 outflow water	2	P4_P402
	Spare actuator	1	P4_P403
	Spare actuator	1	P4_P404
	Unknown	2	P4_UV401

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

Conclusion and Future Work

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

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

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

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

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

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

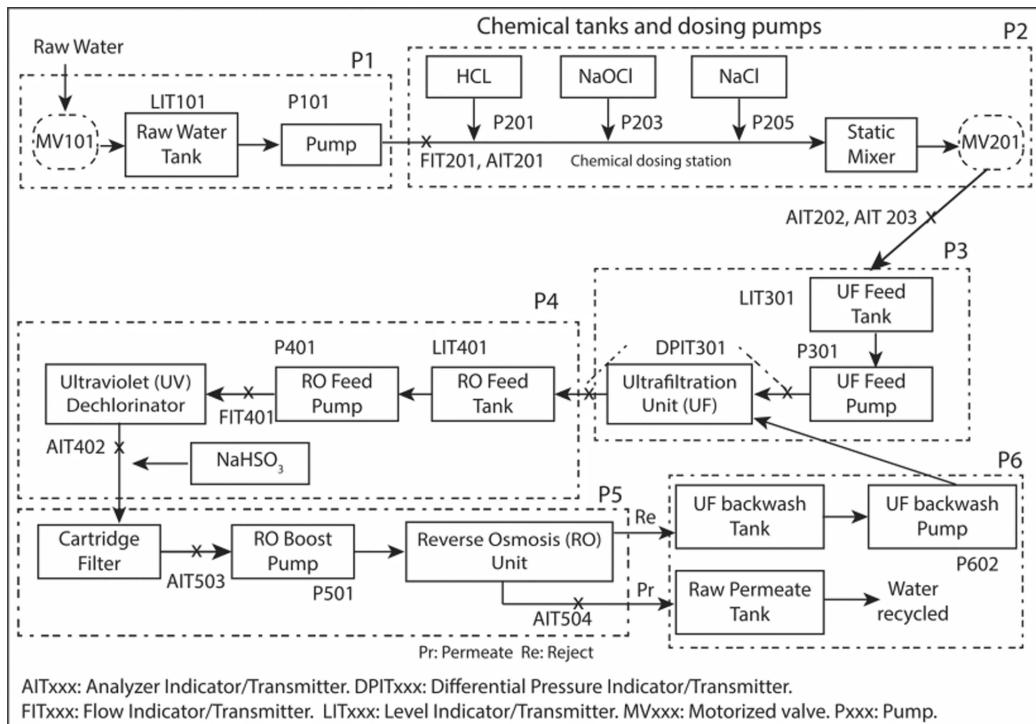


Figure 7.1: iTrust SWaT schema

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

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

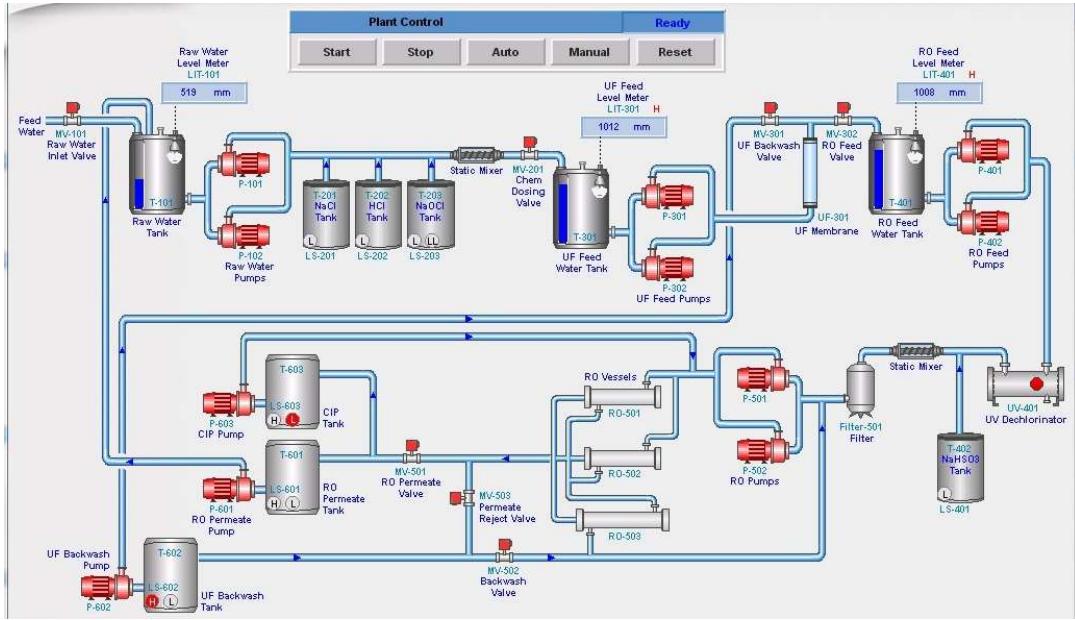


Figure 7.2: iTrust SWaT schema from HMI point of view.

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

Achievements of the approach After comparing the analysis results with the actual system, we can assess the extent to which the proposed framework has achieved its set objectives and identify any limitations encountered.

In terms of the objectives that have been achieved, either fully or partially, we have the following outcomes:

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

- users have the ability to select a specific subsystem of their choice for analysis;
- automatic recognition of registers associated with physical process measurements and actuators (even spare actuators) has been implemented, including their actual operating ranges (*relative setpoints*);
- specific properties of actuators, including their states and the duration of each state, have been identified;
- recognition of which actuators handle incoming or outgoing flows;
- partial recognition of the roles of certain registers within the system has been achieved, such as tanks, valves, pumps, and flow/pressure sensors;
- conditions that cause state changes within the system (*absolute setpoints*) have been recognized;
- data perturbations for specific measurements have been attenuated;
- the Graphical Analysis component now provides a comprehensive view of the system under investigation, enabling the extraction of much more information compared to the corresponding functionality of the Ceccato et al. framework;
- the identification of invariants has been made easier, thanks in part to the implementation of semi-automated analyses;
- the process mining capabilities related to the physical system can extract more data compared to the Ceccato et al. framework;
- the reconstruction of the PLC network and the identification of the industrial protocol used for communications have been achieved.

Limitations of the approach In terms of limitations, we have the following observations:

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

Based on our findings, we can confidently state that the proposed framework has successfully addressed many of the limitations identified in the Ceccato et al. framework. It has provided a greater amount of refined information, aligning with the goals set forth at the beginning of this thesis. However, it is important to acknowledge that certain limitations, particularly related to network data, have restricted the framework from fully realizing its potential. These limitations highlight the need for further advancements in gathering comprehensive and complete network data to unlock the framework's full capabilities.

Future work The framework, along with the thesis sources and analysis files, is openly accessible on the dedicated GitHub repository [75]. 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 used in ICS environments, expanding its compatibility and usability;
- enhance the automatic recognition of sensors and actuators by leveraging advanced machine learning and artificial intelligence techniques. This improvement would involve training models to accurately identify and classify different types of sensors and actuators based on their data patterns and characteristics;
- complete the implementation of network traffic data within the Business Process part of the framework. However, it is crucial to ensure that the network traffic data and physical process data used for analysis are reliable and not compromised by external attackers, as accurate and untampered data is essential for effective analysis;
- implement an automated system for recognizing real system configurations by excluding actuators that do not significantly contribute to changing the system behavior within the analyzed PLC. This automatic recognition system can save time and effort by eliminating the need to manually filter out non-contributing actuators during the analysis process.

Appendix A - Framework Scripts Handbook

Network Analysis

networkAnalysis.py

```
usage: networkAnalysis.py [-h] [-f FILE | -D DIRECTORY] [-
    ↪ s SRCADDR] [-d DSTADDR]

options:
-h, --help
    show this help message and exit
-f FILE, --file FILE
    CSV file with network data
-D DIRECTORY, --directory DIRECTORY
    CSV files with network data
-s SRCADDR, --srcaddr SRCADDR
    Source IP address
-d DSTADDR, --dstaddr DSTADDR
    Destination IP address
```

export_pcap_data.py

```
usage: export_pcap_data.py [-h] [-f FILE | -m MERGEFILES [-
    ↪ MERGEFILES ...] | -d MERGEDIR | -s SINGLEDIR] [-t
    ↪ TIMERANGE TIMERANGE]
```

```
options:  
-h, --help  
    show this help message and exit  
-f FILE, --file FILE  
    single pcap file to include (with path)  
-m MERGEFILES [MERGEFILES ...], --mergefiles MERGEFILES [  
    ↪ MERGEFILES ...]  
    multiple pcap files to include (w/o path)  
-d MERGEDIR, --mergedir MERGEDIR  
    directory containing pcap files to merge  
-s SINGLEDIR, --singledir SINGLEDIR  
    directory containing pcap files for single analysis  
-t TIMERANGE TIMERANGE, --timerange TIMERANGE TIMERANGE  
    time range selection (format YYYY-MM-DD HH:MM:SS)
```

Pre-processing and Preliminary Analysis

mergeDatasets.py

```
usage: mergeDatasets.py [-h] [-g GRANULARITY] [-s SKIPROWS [  
    ↪ ] [-n NROWS] [-d DIRECTORY] [-o OUTPUT] [-p PLCS [  
    ↪ PLCS ...]]  
  
options:  
-h, --help  
    show this help message and exit  
-g GRANULARITY, --granularity GRANULARITY  
    choose granularity in seconds (for slopes)  
-s SKIPROWS, --skiprows SKIPROWS  
    skip seconds from start  
-n NROWS, --nrows NROWS  
    number of seconds to consider  
-d DIRECTORY, --directory DIRECTORY  
    directory containing CSV files  
-o OUTPUT, --output OUTPUT  
    output file  
-p PLCS [PLCS ...], --plcs PLCS [PLCS ...]
```

```
PLCs to include (w/o path)
```

system_info.py

```
usage: system_info.py [-h] [-f FILENAME] [-a ACTUATORS [
    ↪ ACTUATORS ...]] [-s SENSORS [SENSORS ...]]  
  
options:  
-h, --help  
    show this help message and exit  
-f FILENAME, --filename FILENAME  
    name of the input dataset file (CSV format)  
-a ACTUATORS [ACTUATORS ...], --actuators ACTUATORS [  
    ↪ ACTUATORS ...]  
    actuators list  
-s SENSORS [SENSORS ...], --sensors SENSORS [SENSORS ...]  
    sensors list
```

Graphical and Statistical Analysis

runChartSubPlots.py

```
usage: runChartSubPlots.py [-h] [-f FILENAME] [-r
    ↪ REGISTERS [REGISTERS ...]] [-e EXCLUDEREGISTERS [
    ↪ EXCLUDEREGISTERS ...]] [-a ADDREGISTERS [ADDRREGISTERS
    ↪ ...]]  
  
options:  
-h, --help  
    show this help message and exit  
-f FILENAME, --filename FILENAME  
    CSV file to analyze  
-r REGISTERS [REGISTERS ...], --registers REGISTERS [  
    ↪ REGISTERS ...]  
    registers to include  
-e EXCLUDEREGISTERS [EXCLUDEREGISTERS ...], --
```

```
    ↳ excluderegisters EXCLUDEREGISTERS [EXCLUDEREGISTERS
    ↳ ...]
    registers to exclude
-a ADDREGISTERS [ADDREGISTERS ...], --addregisters
    ↳ ADDREGISTERS [ADDREGISTERS ...]
    registers to add
```

Invariant Inference and Analysis

daikonAnalysis.py

```
usage: daikonAnalysis.py [-h] -f FILENAME [-s
    ↳ SIMPLEANALYSIS] [-c CUSTOMANALYSIS] [-u UPPERMARGIN]
    ↳ [-l LOWERMARGIN]

options:
-h, --help
    show this help message and exit
-f FILENAME, --filename FILENAME
    name of the input dataset file (CSV format) w/o path
-s SIMPLEANALYSIS, --simpleanalysis SIMPLEANALYSIS
    simple Daikon analysis on single actuators
-c CUSTOMANALYSIS, --customanalysis CUSTOMANALYSIS
    Daikon analysis based on actuators state changes
-u UPPERMARGIN, --uppermargin UPPERMARGIN
    Upper safety margin (percent)
-l LOWERMARGIN, --lowermargin LOWERMARGIN
    Lower safety margin (percent)
```

runDaikon.py

```
usage: runDaikon.py [-h] [-f FILENAME] [-r REGISTER] [-c
    ↳ CONDITIONS [CONDITIONS ...]]]

options:
-h, --help
```

```
show this help message and exit
-f FILENAME, --filename FILENAME
    name of the input dataset file (CSV format) w/o path
-r REGISTER, --register REGISTER
    sensor to make results directory
-c CONDITIONS [CONDITIONS ...], --conditions CONDITIONS [
    ↪ CONDITIONS ...]
    Daikon invariant conditions
```

Business Process Mining and Analysis

processMining.py

```
usage: processMining.py [-h] [-f FILENAME] [-a ACTUATORS [
    ↪ ACTUATORS ...]] [-s SENSORS [SENSORS ...]] [-t
    ↪ TOLERANCE] [-g GRAPH]

options:
-h, --help
    show this help message and exit
-f FILENAME, --filename FILENAME
    name of the input physical dataset file (CSV format)
-a ACTUATORS [ACTUATORS ...], --actuators ACTUATORS [
    ↪ ACTUATORS ...]
    actuators list
-s SENSORS [SENSORS ...], --sensors SENSORS [SENSORS ...]
    sensors list
-t TOLERANCE, --tolerance TOLERANCE
    tolerance
-g GRAPH, --graph GRAPH
    generate state graph
```


List of Figures

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

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

6.3	Chart of PLC1-2 registers without spare actuators (particular)	133
6.4	Slope calculation anomaly (in the circle)	138
6.5	Activity diagram for PLC1-2	139
6.6	Verifying the conjecture about valve P2_MV201	148
6.7	PLC3 registers	149
6.8	P3_MV301 and P3_MV303 analysis	150
6.9	Activity diagram for PLC2-3	156
6.10	PLC4 registers	161
6.11	P3_MV302 and P4_LIT401 behaviors	162
6.12	Tanks in subsystem PLC3-4 and their correlation.	163
6.13	Activity diagram for PLC3-4	165
7.1	iTrust SWaT schema	170
7.2	iTrust SWaT schema from HMI point of view.	171

List of Tables

2.1	Differences between Information Technology (IT) and Industrial Control Systems (ICSs)	10
2.2	Modbus Function Codes list	24
3.1	Summary table of Ceccato et al. framework limitations	60
4.1	Summary table of proposed framework enhancements	112
5.1	Main IP addresses of the six PLCs and SCADA in the SWaT system	118
5.2	SWaT network traffic data	121
6.1	Properties of the PLC1-2 subsystem	143
6.2	Properties of the PLC2-3 subsystem	158
6.3	Properties of the PLC3-4 subsystem	167
6.4	Summary of the reconstructed composition PLCs from 1 to 4 in the iTrust SWaT System	168

Listings

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

6.5	Conditional Invariants for states 1 and 2 of P1_MV101	136
6.6	Preliminary analysis outcomes for sensors and actuators of PLC2-3	143
6.7	Time duration of the states of actuators of PLC3	144
6.8	P2_MV201 state changes in relation to P3_LIT301	146
6.9	Conditional Invariants for P2_MV201 and P3_P302	151
6.10	Slope calculation of P3_LIT301 for the 850-900 and 970-1000 intervals related to tank levels	154
6.11	Preliminary analysis outcomes for sensors and actuators of PLC3-4	158
6.12	P3_MV302 state changes in relation to P4_LIT401	160
6.13	Daikon manual analysis for P3_MV302 == 2	164

References

- [1] 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.
- [2] R. Rajkumar, L. Lee, I. Sha, and J. A. Stankovic. “Cyber-physical systems:the next computing revolution.” In: 47th Design Automation Conference, DAC 2010, Anaheim, California, USA (July 2010). 2010. DOI: <http://dx.doi.org/10.1145/1837274.1837461>.
- [3] NIST. *ICS defintion*. URL: https://csrc.nist.gov/glossary/term/industrial_control_system (visited on 2023-04-06).
- [4] 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).
- [5] Schneider Electrics. *What is Modbus and how does it work?* URL: <https://www.se.com/us/en/faqs/FA168406/> (visited on 2023-04-13).
- [6] Wikipedia. *DNP3*. URL: <https://en.wikipedia.org/wiki/DNP3> (visited on 2023-05-18).
- [7] *Introduction to EtherNet/IP Technology*. URL: https://scadahacker.com/library/Documents/ICS_Protocols/Intro%20to%20EthernetIP%20Technology.pdf (visited on 2023-04-19).

- [8] OPC Foundation. *Unified Architecture*. URL: <https://opcfoundation.org/about/opc-technologies/opc-ua/> (visited on 2023-06-15).
- [9] Wireshark.org. *S7comm*. URL: <https://wiki.wireshark.org/S7comm> (visited on 2023-06-15).
- [10] 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>.
- [11] 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).
- [12] 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>.
- [13] A. Keliris and M. Maniatakos. “ICSREF: A Framework for Automated Reverse Engineering of Industrial Control Systems Binaries”. In: 26th Annual Network and Distributed System Security Symposium (NDSS) (Feb. 2019). 2019. DOI: <https://doi.org/10.48550/arXiv.1812.03478>.
- [14] *The Daikon dynamic invariant detector*. URL: <http://plse.cs.washington.edu/daikon/> (visited on 2023-04-25).
- [15] iTrust Singapore University of Technology and Design. *Secure Water Treatment*. URL: <https://itrust.sutd.edu.sg/testbeds/secure-water-treatment-swat/> (visited on 2023-05-15).
- [16] Wikipedia. *Stuxnet*. URL: <https://en.wikipedia.org/wiki/Stuxnet> (visited on 2023-05-18).
- [17] Wikipedia. *Industroyer*. URL: <https://en.wikipedia.org/wiki/Industroyer> (visited on 2023-05-18).

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

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

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

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

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

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