



# Provable Adversarial Safety in Cyber-Physical Systems

John Castellanos, Mohamed Maghenem, Alvaro Cárdenas, Ricardo Sanfelice,  
Jianying Zhou

## ► To cite this version:

John Castellanos, Mohamed Maghenem, Alvaro Cárdenas, Ricardo Sanfelice, Jianying Zhou. Provable Adversarial Safety in Cyber-Physical Systems. EuroS&P 2023 - IEEE 8th European Symposium on Security and Privacy (EuroS&P), Jul 2023, Delft, Netherlands. pp.979-1012, 10.1109/EuroSP57164.2023.00062 . hal-04302966

**HAL Id: hal-04302966**

**<https://hal.science/hal-04302966>**

Submitted on 23 Nov 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Provable Adversarial Safety in Cyber-Physical Systems

John H. Castellanos

*CISPA Helmholtz Center for Information Security  
Saarbrücken, Germany  
john.castellanos@cispa.de*

Mohamed Maghenem

*University of Grenoble Alpes, CNRS  
Grenoble INP, France  
mohamed.maghenem@gipsa-lab.grenoble-inp.fr*

Alvaro A. Cárdenas  
*University of California  
Santa Cruz, USA  
alacarde@ucsc.edu*

Ricardo G. Sanfelice  
*University of California  
Santa Cruz, USA  
ricardo@ucsc.edu*

Jianying Zhou  
*Singapore University of Technology and Design  
Singapore  
jianying\_zhou@sutd.edu.sg*

**Abstract**—Most proposals for securing control systems are heuristic in nature, and while they increase the protection of their target, the security guarantees they provide are unclear. This paper proposes a new way of modeling the security guarantees of a Cyber-Physical System (CPS) against arbitrary false command attacks. As our main case study, we use the most popular testbed for control systems security. We first propose a detailed formal model of this testbed and then show how the original configuration is vulnerable to a single-actuator attack. We then propose modifications to the control system and prove that our modified system is secure against arbitrary, single-actuator attacks.

## 1. Introduction

In the past decades, we have seen several confirmed attacks on industrial control systems, including a sewage control system in Australia [57], a nuclear enrichment facility in Iran [69], the power grid of Ukraine [13], a steel mill in Germany [39], a paper mill in Louisiana [6], oil systems in the Middle-East [35], and a water utility in Florida [45]. In all these cases, an attacker partially compromised a control system and then sent malicious control commands to the physical process, causing accidents and damages.

Researchers have suggested various defense strategies; however, two main challenges remain largely unaddressed. First, most security efforts for industrial control security are heuristic in nature, and they do not provide provable security assertions about the system’s safety. Second, implementing and evaluating security proposals for industrial control systems is generally restricted to either simulations or toy physical systems, given the difficulty of getting access to real-world operational industrial plants.

This paper addresses these two limitations by proposing a new method to formally prove safety properties against an attacker that has partially compromised the system and implementing and testing our formal model in a real-world operating plant. In particular, we implement our proposal in the most popular industrial process for CPS security research [15]: the Secure Water Treatment testbed (SWaT).

First, this paper proposes the most comprehensive formal mathematical model of SWaT. As we show in

Table 1, there is no previous model that is larger than ours (our representation of SWaT captures in a single model as many or more SWaT elements than previous work). We have also released our model as open-source software (see itemized contributions below). In addition, while several papers in top security conferences have used SWaT before [3], [9], [12], [21], [63], none of these previous efforts attempted to provide a comprehensive formal model of the system nor enable proofs of security assertions.

We then formally prove the following sequence of results: (1) We prove that without attacks, SWaT is safe. (2) We show that if the attacker can compromise a single actuator in the system, then the original design of SWaT is unsafe (some attacks drive the system to unsafe regions). (3) We propose modifications to the Programmable Logic Controllers (PLCs) operating SWaT and then formally prove that SWaT will remain safe under single actuator attacks with these changes. (4) We show that if the attacker can compromise two or more actuators, SWaT is unsafe, regardless of the logic in the PLCs.

To prove these security assertions we extend the theory of barrier function certificates for hybrid systems and adapt them to analyze the safety of a system under attack. Our contributions include extensions to the theory of barrier functions by introducing the concept of uniform safety in the presence of arbitrary exogenous signals.

Furthermore, we propose a new adversary model that does not make parametric assumptions about the attacker’s tactics. Our adversary model only needs to know the number of actuators an attacker has under its control but does not need to know the tactics. In other words, the attacker can launch a square wave attack of any frequency, a delay attack, an inversion attack, etc. Any arbitrary attack signal is considered in our model, making our security proofs robust to new unanticipated attack tactics. This is contrary to most work on industrial security which makes parametric assumptions about the attack strategies (e.g., scaling attacks [60], bias attacks [12], [14], delay attacks [60], or random attacks [20], [67]). We note that actuators can be compromised individually by targeting the Remote Input/Output (RIO) computer interfacing the actuator with the PLC [64].

Finally, we evaluate and test our methods and proofs in a real-world system. We also show evidence that our model is accurate by comparing our model and traces from the real-world system.

In summary, our contributions include:

- 1) We propose a framework to combine the *static analysis of PLC code* to create **control invariants**, with *traces of the physical behavior of the system* to create **physical invariants**.
- 2) Using the two approaches above, we introduce the most comprehensive formal model of a popular [15] testbed for CPS security.
- 3) We formally prove several security assertions about SWaT, including (1) safety without attacks, (2) unsafe if an attacker compromises a single actuator, (3) safe with our PLC modifications (under a single actuator attack), (4) unsafe with two or more actuators under attack.
- 4) We extend the theory of barrier function certificates for hybrid systems by introducing a new concept for proving safety in adversarial conditions. We then find sufficient conditions to prove safety under our adversary model.
- 5) Our adversary model considers attackers that can launch arbitrary control commands. Our adversary model is more realistic and powerful than previous adversaries considered in industrial control systems [12], [14], [20], [21], [60], [67].
- 6) We implement and validate our approach in a real-world system. We release our model as open-source software<sup>1</sup> in the hopes that other researchers working with SWaT can use it.

## 1.1. Organization of the paper

The rest of the paper is organized as follows: In Section 2 we discuss related work. Then in Section 4 we introduce the mathematical tools to model a cyber-physical system, our adversary model (Section 4.1), and a theorem to show how to prove the safety of control systems under attacks. In Section 5 we introduce the physical model of SWaT. We also prove that this system is safe without attacks for the first time. We evaluate the security of our water system in Section 6 and show how the system is safe under some attacks and unsafe under other attacks. We then propose changes to the control logic and to the physical parameters of the system to make the system safe, irrespective of any individual control signal being compromised. We then finalize our analysis with experimental results validating our theoretical model in Section 7.

## 2. Related Work

### 2.1. Attacks to CPS

There are various ways in which attackers can take over sensors or actuators in control systems. In addition to compromising devices in the “classical” way (e.g., a

software exploit), attackers can also compromise sensor or control signals with novel physical attacks. Analog sensor security [68], [23], [62], [58], [8] focuses on how physical interference can affect the reported sensor readings back to the control system. By adding a physical signal (electromagnetic, sound, heat, etc.), attackers can affect sensor readings. In some cases, these types of attacks can also manipulate actuators directly [55], [16].

When attackers compromise a control signal in a CPS (e.g., acceleration in a vehicle), they inject a time series to the physical system  $a(t)$ , where  $t$  denotes time. Researchers tend to parameterize the adversary tactics to fixed strategies, and they assume that the attacker simply replaces the non-compromised signal  $u(t)$  with a parameterized version of it. Examples include scaling attacks  $a(t) = \alpha u(t)$  [60], bias attacks  $a(t) = u(t) + b$  [12], [14], delay attacks  $a(t) = u(t - d)$  [60], and random attacks (where  $a(t)$  is a random value at each time) [67], [20].

As history has taught us, limiting the attacker to follow specific attacks will not guarantee security. Over the years, system after system has been defeated by adversaries that break the assumptions of the model. A recent high-profile example is the case of attacks against key handshakes in WPA2 [65] which were proven secure [30] under a model that did not capture key installation. To increase the confidence of a security proof, adversary models used in formal proofs tend to be as general as possible [17], [37], [27], e.g., by assuming adversaries to be **any** polynomial time algorithm [27], [37] (without parameterizing the specific attacker algorithm used to crack the system). Our goal in this paper is to study attacks without assuming a priori a fixed attack tactic.

### 2.2. Securing CPS

Defenses for control systems can be reactive or proactive. Reactive security focuses on detecting attacks [12], [14], [21], [63], [29], [54], [51] and sometimes responding to attacks [70], [51], [20], [19]. Attack response usually focuses on first identifying the malicious sensor or actuator signals and then eliminating them.

Proactive security proposals, on the other hand, focus on secure design: for example, designing a control algorithm so that the system is more resilient to attacks or designing actuators so what the attacker can do is limited [36], [24]. This paper focuses on proactive security: *we want to evaluate the system offline to understand the impact of attacks and if possible, redesign the controller to minimize the negative impacts of any future attack.*

In summary, the **scope** of this paper is to study (offline) the safety of an industrial design, identify the stress points where an attacker can break the system, and propose improvements to make the system more secure and resilient to attacks. We approach this problem by providing provable security assertions.

### 2.3. Formal Verification of CPS

Embedded control systems monitor and control a variety of safety-critical problems. To formally guarantee the safe operation of CPS, we need formal models and

1. <https://gitlab.com/jhcastel/provable-adversarial-safety-in-cps>

rigorous verification approaches. In this subsection we summarize different ways to formally verify a CPS and introduce how the tool we use in the paper (barrier certificates) compares to alternatives.

In most classical problems in computer science, model checking focuses on discrete dynamics; however, the unique challenge of verifying CPS properties arises due to their continuous dynamics. Continuous dynamics creates new problems as we need to consider the evolution of physical states that follow trajectories defined by differential equations.

There are three main approaches for the verification of CPS: (1) set-based reachability analysis, (2) abstraction-based verification, and (3) logic-based verification [18].

Reachability analysis attempts to find the set of reachable states of a CPS as it evolves over time. The goal is to check if the reachable set and the unsafe set are disjoint. There are several tools for reachability calculations; for example, Flow\* [11] uses a flowpipe construction scheme to verify time-bounded reachability. Reachability analysis has also been applied in an adversarial setting [70], [36], [24], [59]; however, these security efforts for applying reachability in adversarial conditions have only considered linear systems, while real-world CPS have more complex physical behaviors (like hybrid dynamics). In this paper, we go beyond linear systems and apply our methods to a real-world control system modeled by hybrid equations.

Abstraction-based verification attempts to address the scalability problems of reachability analysis. To scale up discrete model checking, abstraction-based verification replaces the actual system with a simpler, abstract system in which model checking is easier to perform [31], [4]. The drawback of these abstractions is that any verification result can only be related back to the original system if the property in question survives the abstraction process [18].

Finally, logic-based verification provides a witness to verify a continuous system respects the desired property. There are two main methods for logic-based verification: (a) differential invariants [46], [47], and (b) barrier certificates [48], [5]. Differential invariants are based on Lie derivatives and Lie groups, while barrier certificates are based on Lyapunov’s criterion for stability [18]. Differential invariants are the most general representation for logic-based verification, but if you find a barrier certificate, logic-based verification can be proven directly (analytically) and without computational support (e.g., requiring software tools). In addition, software tools require the discretization of continuous states, while analytical barrier certificate proofs do not require this approximation [33].

Barrier certificates are a way to separate good and bad states and to show that this separation (barrier) is impenetrable by the continuous system dynamics [42]. The importance of barrier certificates comes from the fact that they reduce a reachability question (can we ever reach an unsafe state) by a simple check on the directional derivative of the barrier certificate along the differential equation of the system [18]. In contrast with other tools, barrier certificates provide formal analytical guarantees without extra assumptions or without the need to rely on computational tools. For example, reachability tools like Flow\* [11] use a flowpipe construction scheme to verify time-bounded reachability. Flowpipe construction

methods are often easily used since users only need to specify the flow pipe stepsize, approximation order, and the bounded time horizon. On the other hand, barrier certificates require an expert to find a barrier function, but a barrier certificate can be used to prove time-unbounded reachability. In short, barrier certificates, like the ones we obtain in this work, make safety proofs harder, but their safety guarantees are stronger than time-bounded reachability alternatives, proofs relying on abstractions, or proofs relying on bounded computational checks.

In this paper, we develop the theory of barrier certificates for CPS verification of safety under actuation attacks. We also contribute to the literature on safety by defining the new concept of *uniform safety*, which is required when an attacker is the source of uncertainty in the system. In addition, we prove new theorems that show how to check if a barrier certificate satisfies uniform safety.

In this paper we focus on formal verification, rather than testing. There have been other efforts to use testing for SWaT. For example, HyChecker [38] combines random sampling with symbolic execution on hybrid systems to perform probabilistic security testing. HyChecker is not a formal verification approach, so it cannot guarantee safety properties in the system. Barrier certificates can provide these safety guarantees [42].

## 2.4. Previous Models of SWaT

Previous works	Modeling approach	AL	SWaT Modeled	Stages	U
Adepu [1], Feng [21]	Black-Box	C	Tank level, chemical flows		D
Ahmed [3], [2]	Black-Box	O	Sensor signals		D
Castellanos [10]	Black-Box	F	Tank level (Stage 1)		D
Chen [12]	Black-Box	F	Tank level (Stage 1, 3, 4)		D
Lin [40]	Black-Box	F	Tank level (Stage 1, 3, 4)		D
Urbina [63]	Black-Box	F	Chemical dosing (Stage 2)		D
HyChecker [38]	Probabilistic hybrid model	C	Backwash tank		T
This work	First-Principles	F	Tank level (Stage 1, 3, 4)		S

TABLE 1: SWaT models studied in previous security conferences. Abstraction level(AL): (O) Orthogonal model, (C) Coarse model, (F) Fine-grained model. Use case(U): (D) Attack Detection, (T) Security Testing, (S) Formal Proofs of Adversarial Safety.

SWaT is a real-world water system that has been widely used in security and formal method conferences [1], [2], [3], [9], [10], [12], [21], [38], [40], [63].

Some of these efforts have attempted to model SWaT. Table 1 shows how our model of SWaT compares to other models of SWaT. Our paper models the interactions of three interconnected stages in a single fine-grained mathematical model and therefore is the model capturing most of the physics in SWaT (other papers simply

model a smaller subset of SWaT). Only two other related works [12], [40] have provided a similar comprehensive model of SWaT. Compared to them, our model is explicitly given by mathematical equations derived from first principles, while previous work attempting to capture the same level of complexity has resorted to black-box machine learning models that do not provide any guarantees about the accuracy of the model, or the explanation of the dynamics (therefore these previous models cannot be used for mathematical proofs of security).

Coarse models [22], [1] do not provide accuracy for the time series in the testbed and, therefore cannot model the precise effects of adversaries. And others [2], [3] are orthogonal models that do not attempt to model SWaT but instead attempt to fingerprint the innate noise of sensors.

In short, our work is the most complete model of SWaT available in the literature, and furthermore, it is a formal mathematical model that can be used in formal proofs of security, and it is based on first principles, so it can be used to explain the interactions between various components.

### 3. Background

**Notation.** Let  $\mathbb{R}_{\geq 0} := [0, \infty)$  and  $\mathbb{N} := \{0, 1, \dots, \infty\}$ . Given two vectors  $x$  and  $y$  of the same dimension,  $m_x$  denotes the dimension of  $x$ ,  $x^\top$  denotes the transpose of  $x$ ,  $|x|$  denotes the Euclidean norm of  $x$ , and  $\langle x, y \rangle = x^\top y$  denotes the scalar product of  $x$  and  $y$ . Given a nonempty set  $K \subset \mathbb{R}^{m_x}$ ,  $\text{int}(K)$  denotes the interior of  $K$ ,  $\partial K$  denotes its boundary,  $\text{cl}(K)$  denotes its closure, and  $U(K)$  denotes an open neighborhood of the set  $K$ . For a nonempty set  $O \subset \mathbb{R}^{m_x}$ ,  $K \setminus O$  denotes the subset of elements of  $K$  that are not in  $O$ . For a differentiable map  $x \mapsto B(x) \in \mathbb{R}$ ,  $\nabla B$  denotes the gradient of  $B$  with respect to  $x$ . Finally, by  $\dot{x}$ , we denote the time derivative of the state  $x$ , while by  $x^+$  we denote the value of the state after an instantaneous jump.

#### 3.1. Hybrid-model Approach

In this section, we propose a new hybrid-model approach to analyze the safety of cyber-physical systems (CPS) in the presence of attacks. Hybrid systems are models that enrich computing models with analog models of physics; as a result, they contain digital models of computing (such as automata or programs) as well as analog elements (such as differential equations) integrated in a way that allows us to reason about the effect of physics on computing and vice versa [18].

Formally, a hybrid equation is composed of a differential equation with a constraint, which models the *flow* or the continuous evolution of the system (e.g., level of water in a tank), and a difference equation with a constraint, modeling the *jumps* or discrete events (e.g., a change in the status of an actuator from ON to OFF). The strength of the hybrid equations formalism relies on the compactness of the representation and the possibility of separately using or extending the existing tools developed for continuous and discrete-time systems. Following [26], a hybrid dynamical system  $\mathcal{H} = (C, F, D, G)$  as in (3) with the state variable

$x \in X \subset \mathbb{R}^{m_x}$ , the flow set  $C \subset X$ , the jump set  $D \subset X$ , the flow and jump maps  $F : X \rightarrow X$  and  $G : X \rightarrow X$ , respectively.

A hybrid arc  $\phi$  is defined on a hybrid time domain denoted  $\text{dom } \phi \subset \mathbb{R}_{\geq 0} \times \mathbb{N}$ . The hybrid arc  $\phi$  is parametrized by an ordinary time variable  $t \in \mathbb{R}_{\geq 0}$  and a discrete jump variable  $j \in \mathbb{N}$ . Its domain of definition  $\text{dom } \phi$  is such that for each  $(T, J) \in \text{dom } \phi$ ,  $\text{dom } \phi \cap ([0, T] \times \{0, 1, \dots, J\}) = \bigcup_{j=0}^J ([t_j, t_{j+1}] \times \{j\})$  for a sequence  $\{t_j\}_{j=0}^{J+1}$ , such that  $t_{j+1} \geq t_j$ ,  $t_0 = 0$ , and  $t_{j+1} = T$ .

We define the concept of a *solution*  $x$  to a hybrid equation  $\mathcal{H} := (C, F, D, G)$ .

**Definition 1 (Concept of solutions to  $\mathcal{H}$ ).** A hybrid arc  $x : \text{dom } x \rightarrow X$  is a *solution* to  $\mathcal{H}$  if

(S0)  $x(0, 0) \in \text{cl}(C) \cup D$ ;

(S1) for all  $j \in \mathbb{N}$  such that  $I^j := \{t : (t, j) \in \text{dom } x\}$  has nonempty interior,  $t \mapsto x(t, j)$  is locally absolutely continuous and

$$\begin{aligned} x(t, j) &\in C && \text{for all } t \in \text{int}(I^j), \\ \dot{x}(t, j) &= F(x(t, j)) && \text{for a.a. } t \in I^j; \end{aligned} \quad (1)$$

(S2) for all  $(t, j) \in \text{dom } x$  such that  $(t, j+1) \in \text{dom } x$ ,

$$x(t, j) \in D, \quad x(t, j+1) = G(x(t, j)). \quad (2)$$

•

#### 3.2. Modeling CPS using Hybrid Equations

In CPS, we identify two types of state variables. The physical variables; e.g., the water levels, which change continuously with respect to time and take values from a dense set; e.g.  $\mathbb{R}_{\geq 0}$ . On the other hand, the logic variables, e.g., the state of motor valves or pumps, which change through a code executed at PLCs. Those discrete variables take values from discrete sets, e.g., ON, OFF, or in transition. This heterogeneous combination of variables requires dynamical models combining continuous and discrete variables [53]. Hence, hybrid system models [26] are a natural framework for studying cyber-physical systems.

### 4. Proving Safety Under Attacks

#### 4.1. Adversary Model

In the presence of actuator attacks (see Fig. 1), the attacker can falsify the actuation given to the system, either by compromising the control signal sent by the controller (right side in Fig. 1) or by compromising the actuator directly with a digital or a transduction attack [16], [55] (left side in Fig. 1). Throughout the paper, we assume the attacker can compromise one actuator and change its control action. The attacker can achieve this partial compromise by exploiting memory vulnerabilities or resource access control vulnerabilities (based on the ICS vulnerabilities categorization [61]). In this paper we focus on *post-exploitation* rather than on the specific method the attacker used to get access into the system. Our goal is to

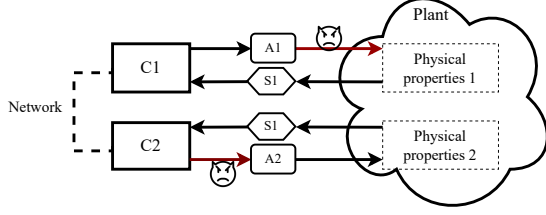


Figure 1: Simplified diagram of a CPS under actuator attacks with (C) Controllers, (S) Sensors, and (A) Actuators. Attackers can directly affect the plant via transduction attack (left side) or controller signal manipulation (right side).

understand if a partially compromised system can remain safe.

One of our goals is to analyze security under a wide variety of attacks. In the last decade, a variety of adversary tactics have been proposed in the literature. Most of them assume that the control signal of the attacker is constrained to few parametric models. For example, a scaling attack [60] takes a compromised signal and scales it with a constant, a bias attack [12], [14] takes a compromised signal and adds a constant bias, abrupt-attacks take the maximum possible value the compromised signal can have [12], [14], [21], delay attacks take a compromised signal and delay it in time [60], and random attacks replace the compromised signal by a signal chosen from a random probability distribution [67], [20]. While all of the examples presented so far are from cybersecurity conferences, the literature in control systems has very similar attack models with delay attacks [43], [32], or scaling attacks [28]. In this paper we do not place any constraints on the control signal sent by the attacker. In our case, by proposing a more general adversary model, we are not confined to existing predefined attacks.

## 4.2. Analyzing Safety in a CPS Under Attack

**4.2.1. Safety Without Attacks.** Now, we turn our attention to the concept of safety which is the property we want to analyze when the system is under attack. Intuitively safety means that the physical process will not cause harm to humans, the environment, or damage the equipment. This is best characterized by keeping a set of state variables inside a boundary (the safe set); if these variables (e.g., the level of water in a tank) go outside the safe set, then it means the system reaches an unsafe condition (i.e., the variables are in an unsafe set). To formalize this notion, we consider a hybrid system given by

$$\mathcal{H} : \begin{cases} \dot{x} = F(x) & x \in C \\ x^+ = G(x) & x \in D, \end{cases} \quad (3)$$

and we let two sets  $X_o \subset \text{cl}(C) \cup D \subset X$  and  $X_u \subset X \setminus X_o$ . The set  $X_o$  represents the set of initial conditions and the set  $X_u$  represents the unsafe set.

**Definition 2 (Safety [48]).**  $\mathcal{H}$  is said to be safe with respect to  $(X_o, X_u)$  iff solutions starting at  $X_o$  never reach  $X_u$ . •

One of the main analytical tools to study safety in hybrid systems is the concept of barrier functions. A barrier function is as a scalar function of the state of the system with a given sign on the set of initial conditions  $X_o$  and the opposite sign on the unsafe set  $X_u$ .

**Definition 3 (Barrier function candidate [41]).** A function  $B : X \rightarrow \mathbb{R}$  is a barrier function candidate with respect to  $(X_o, X_u)$  iff

$$\begin{aligned} B(x) &> 0 & \forall x \in X_u \cap (\text{cl}(C) \cup D) \\ B(x) &\leq 0 & \forall x \in X_o. \end{aligned} \quad (4)$$

Note that the barrier function candidate  $B$  in Definition 3 defines the zero-sublevel set

$$K_e := \{x \in X : B(x) \leq 0\}. \quad (5)$$

Notice that  $X_o \subset K_e$  and  $X_u \cap (C \cup D) \cap K_e = \emptyset$ . Hence, safety is guaranteed provided that the barrier function candidate remains nonpositive when evaluated along the solutions starting from the initial set  $K_e$ ; namely, the set  $K_e$  is forward pre-invariant. Our previous results [anonymized] identified the following sufficient conditions to certify forward pre-invariance of  $K_e$ , which in turn imply safety of  $\mathcal{H}$ :

$$\langle \nabla B(x), F(x) \rangle \leq 0 \quad \forall x \in (U(\partial K_e) \setminus K_e) \cap C, \quad (6)$$

$$B(G(x)) \leq 0 \quad \forall x \in D \cap K_e, \quad (7)$$

$$G(x) \subset C \cup D \quad \forall x \in D \cap K_e. \quad (8)$$

To show the intuition of this approach, we illustrate how to check conditions (6)-(8) to prove the thermostat system is safe (without attacks).

**4.2.2. Safety With Attacks.** As our first contribution, we now adapt our previous results to reason about safety under attacks. First, we consider a hybrid system under general attacks

$$\mathcal{H}_u : \begin{cases} \dot{x} = F(x, u) & (x, u) \in C \\ x^+ = G(x, u) & (x, u) \in D. \end{cases} \quad (9)$$

Where the attack  $u$  can affect the physical states, as well as the discrete software logic. To analyze safety in the presence of attacks, we introduce a new concept we call uniform safety.

**Definition 4 (Uniform Safety).** System  $\mathcal{H}_u$  in (9) is said to be safe with respect to  $(X_o, X_u)$  uniformly in  $u \in \mathcal{U}$  iff, for each solution pair  $(x, u)$  to  $\mathcal{H}_u$  such that  $x(0, 0) \in X_o$ , the solution  $x$  never reaches the set  $X_u$ . •

Another contribution in this paper is the derivation of new sufficient conditions to certify uniform safety of  $\mathcal{H}_u$  in the presence of attacks:

$$\langle \nabla B(x), F(x, u) \rangle \leq 0 \quad \forall (x, u) \in (U(\partial K_u) \setminus K_u) \cap C, \quad (10)$$

$$B(G(x, u)) \leq 0 \quad \forall (x, u) \in D \cap K_u, \quad (11)$$

$$G(x, u) \subset C \cup D \quad \forall (x, u) \in D \cap K_u, \quad (12)$$

where  $K_u := K_e \times \mathcal{U}$ ,  $(U(\partial K_u) \setminus K_u) = (U(\partial K_e) \setminus K_e) \times \mathcal{U}$ , and  $C_u := \{x \in X : \exists u \in \mathcal{U} : (x, u) \in C\}$ .

**Definition 5 (Barrier function certificate for safety).** A  $C^1$  barrier function candidate with respect to  $(X_o, X_u)$  becomes a **barrier function certificate for safety** with respect to  $(X_o, X_u)$  if (6)-(8) are satisfied. •

In this paper, we extend our previous results in Theorem 5 so that we are able to prove safety in the presence of attacks.

**Theorem 1.** Given a hybrid system  $\mathcal{H}_u = (C, F, D, G)$  as in (9), suppose that  $F$  is continuous and that there exists a  $C^1$  barrier function candidate  $B$  with respect to  $(X_o, X_u)$  as in (4). The hybrid system  $\mathcal{H}_u$  is safe with respect to  $(X_o, X_u)$  uniformly in  $u \in \mathcal{U}$  if (11) and (12) hold and

$$\langle \nabla B(x), F(x, u) \rangle \leq 0 \quad \forall (x, u) \in (U(\partial K_u) \setminus K_u) \cap C : F(x, u) \in T_{C_u}(x), \quad (13)$$

where  $K_u := K_e \times \mathcal{U}$ ,  $(U(\partial K_u) \setminus K_u) = (U(\partial K_e) \setminus K_e) \times \mathcal{U}$ , and  $C_u := \{x \in X : \exists u \in \mathcal{U} : (x, u) \in C\}$ . □

The proof of Theorem 1 can be found in Appendix C.

In the following example, we illustrate how to use conditions (10)-(12) to prove the uniform safety of a mobile robot under attacks affecting the angular-velocity actuator.

**Example 1 (Proving Safety of a CPS Under Attack).**

Consider a robotic vehicle modeled by the kinematics equation

$$\mathcal{H}_{vu} : \begin{cases} \dot{x} = v \cos(\theta) \\ \dot{y} = v \sin(\theta) \\ \dot{\theta} = u \end{cases} \quad (x, y, \theta, v, u) \in \mathbb{R}^5, \quad (14)$$

where  $v$  and  $u$  are the forward and angular velocities, respectively. The first two elements of the state vector  $[x \ y \ \theta]^\top$  correspond to the Cartesian coordinates of a point on the robot with respect to a fixed reference frame, and  $\theta$  denotes the robot's orientation with respect to the same frame (see Fig. 2).

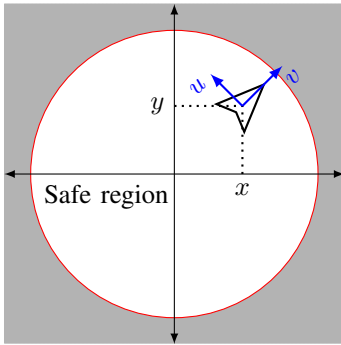


Figure 2: Robotic Vehicle and Safe Region.

The safety requirement consists in maintaining the distance between the vehicle's position and the origin within a given range. This models the case where the

operator of the vehicle is at the origin and the wireless signal for operating the robot will only extend up to a radius  $R$ . If the vehicle wanders outside the safety circle, the operator will lose control of the vehicle. Hence, we assume that

$$\begin{aligned} X_u &:= \{(x, y, \theta) \in \mathbb{R}^3 : |(x, y)|^2 > 1\} \\ X_o &:= \{(x, y, \theta) \in \mathbb{R}^3 : |(x, y)|^2 \leq 1, \\ &\quad x \cos \theta + y \sin \theta = 0\}. \end{aligned}$$

We now transform the coordinates of the robot from the global to the local coordinate frame; that is, we define

$$\begin{bmatrix} x_l \\ y_l \end{bmatrix} := \begin{bmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix}. \quad (15)$$

In these new coordinates, the kinematics equation becomes

$$\mathcal{H}_u : \begin{cases} \dot{\theta} = u \\ \dot{x}_l = uy_l - K(x, y, \theta) \\ \dot{y}_l = -ux_l \end{cases} \quad (x, y, \theta, u) \in \mathbb{R}^4.$$

If the original control is  $v := -x_l$  and  $u = 0$  we show that  $\mathcal{H}_u$  is safe with respect to  $(X_o, X_u)$  when there are no attacks. Indeed, for all  $(x, y, \theta) \in X_o$ ,  $F(x, y, \theta, 0) = 0$ . Hence, the solutions starting from  $X_o$  remain in  $X_o$ . However, as seen in Fig. 3(b), the system is not safe when  $u$  is under attack.

We now redesign  $v$  by choosing  $K(x, y, \theta) := x_l$ . We prove that  $\mathcal{H}_u$  is safe, even under arbitrary attacks in the angular velocity. Indeed, we consider the barrier function candidate

$$B(x, y, \theta) := |(x, y)|^2 - 1 = |(x_l, y_l)|^2 - 1.$$

We now prove that  $\mathcal{H}_u$  is uniformly safe with respect to  $(X_o, X_u)$  by verifying (10)-(12). Indeed, note that (11) and (12) hold trivially since  $\mathcal{H}_u$  is a continuous-time system. Finally, to verify (10), we note that

$$\langle \nabla B(x, y, \theta), F(x, y, \theta, u) \rangle \leq -x_l^2 \leq 0 \quad \forall (x, u) \in \mathbb{R}^4. \quad \square$$

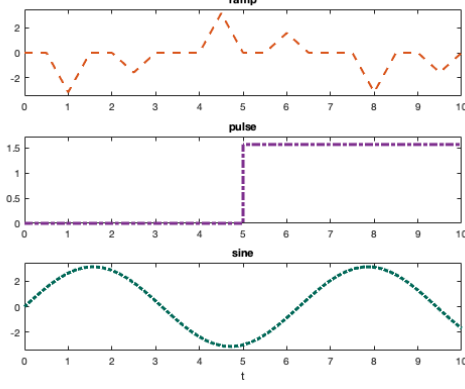
In this example, we have shown how to formally prove the safety of a vehicle when the adversary can use any arbitrary tactic to attack the angular velocity. Notice that “arbitrary” is a key concept to guarantee security. In Fig. 3(c) we see that the vehicle remains safe under three different attacks in Fig. 3(a), but how do we know the system will be safe to another attack we didn't simulate? Theorem 1 (its associated proof in Appendix C) guarantees that the system will remain safe even for attacks we have not simulated.

We now turn the methodology introduced in this section to study SWaT.

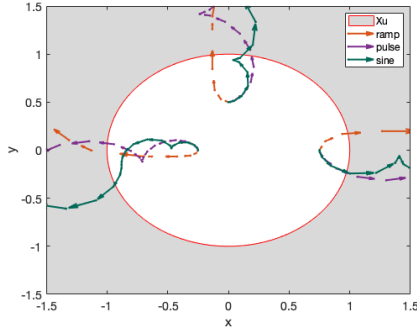
## 5. A Formal Model of SWaT

In this section, we introduce a formal mathematical model of SWaT. Although SWaT has been studied extensively in security conferences [1], [2], [3], [9], [10], [12], [21], [40], [63], to the best of our knowledge, we are the first to derive and share all the equations modeling the

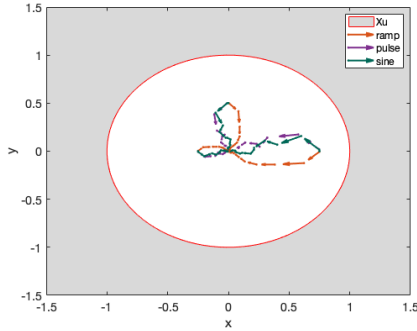




(a) Attack tactics for the angular velocity.



(b) With the first design, the attacker can drive the vehicle outside the safe region.



(c) With the second design, the system is safe under any attack on the angular velocity. Here we see responses from a ramp attack, a pulse attack and a sine attack. But the system will remain safe for any arbitrary tactic from the attacker.

Figure 3: Robotic Vehicle Example.

system. We hope the new and open model in this paper will help future researchers working with SWaT.

SWaT is illustrated in Fig. 4. The control of the flow of water in the process has three stages, and each stage uses a tank to store water with different properties. Stage 1 stores raw water or pre-processing liquid, Stage 2 treats water with chemicals, and Stage 3 stores the water after filtration. The level of water in the three tanks, denoted by  $L1$ ,  $L2$ , and  $L3$ , respectively, has to remain within a given range.

The following components are associated with each stage:

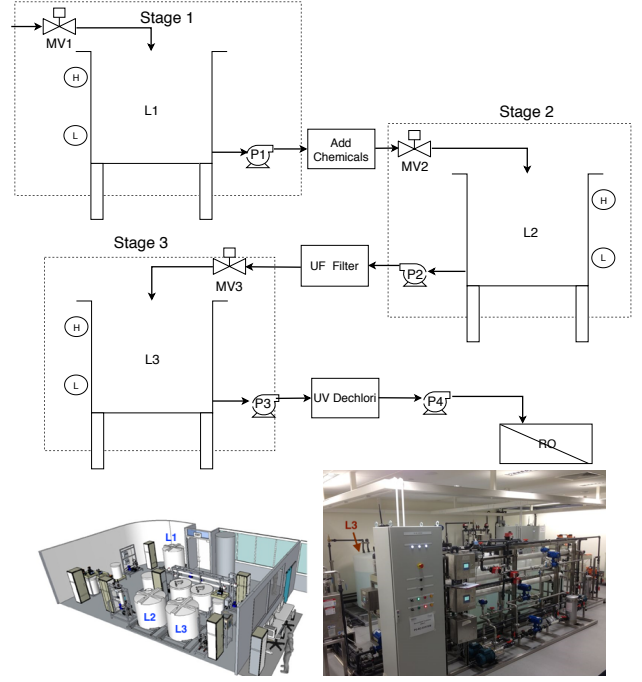


Figure 4: SWaT overview; different PLCs manage each stage.

- Motor valves  $MV1$ ,  $MV2$ , and  $MV3$  feed tanks in each stage. Furthermore, each motor valve has four operation modes: ON ( $\equiv 1$ ), OFF ( $\equiv 0$ ), a transition from ON to OFF denoted  $T\downarrow$  ( $\equiv 3$ ), and a transition from OFF to ON denoted  $T\uparrow$  ( $\equiv 2$ ). Namely,  $(MV1, MV2, MV3) \in \{0, 1, 2, 3\} \times \{0, 1, 2, 3\} \times \{0, 1, 2, 3\}$ .
- Pumps  $P1$ ,  $P2$ , and  $P3$  drain the water from tanks to the next stage.  $P1$  between Stage 1 and Stage 2,  $P2$  between Stage 2 and Stage 3, and  $P3$  between Stage 3 and the final destination. Furthermore, each pump has two operation modes: it allows the water to flow when it is ON ( $\equiv 1$ ), and it blocks the water flow when it is OFF ( $\equiv 0$ ). Hence,  $(P1, P2, P3) \in \{0, 1\} \times \{0, 1\} \times \{0, 1\}$ .
- PLCs  $C1$ ,  $C2$ , and  $C3$  control the water levels ( $L1, L2, L3$ ) in Stages 1, 2, and 3, respectively, by sending commands to  $(MV1, P1)$ ,  $(MV2, P2)$ , and  $(MV3, P3)$ , respectively.

Strictly speaking, for each  $i \in \{1, 2, 3\}$ , the variables  $(P_i, MV_i)$  denote the control signals delivered by the controllers  $C_i$  to the  $i$ -th pump and the  $i$ -th motor valve respectively, as illustrated in Fig. 5.

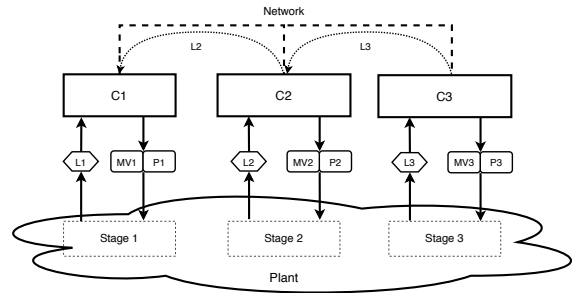


Figure 5: Computer Network of SWaT. PLCs read  $L1-L3$  using sensors, and control the plant through  $MV$  and  $P$ .



## 5.1. Control Invariants

**Stage 3.** PLC C3 uses MV3 to actuate the motor valve and P3 to actuate the pump in Stage 3. The decision to activate or deactivate the motor valve and the pump is based on the information received from a local sensor measuring the water level L3 in the third tank.

Motor valves have four states: in addition to ON and OFF, the additional states are transition steps when migrating from 1 to 0 and vice versa. Moreover, for each  $i \in \{1, 2, 3\}$ , the transition from  $MVi = 3$  to  $MVi = 2$  cannot happen instantaneously, in the sense that the system waits  $T_i$  seconds since  $MVi = 3$  to switch  $MVi$  to 0 and then to 2. The same logic applies when transitioning from  $MVi = 2$  to  $MVi = 3$ .

Each transition mode lasts for  $T_3 > 0$  seconds. In summary, the control logic to actuate on the motor valve for stage 3 is

$$MV3 := \begin{cases} 1 & \text{if } (\tau_3 \geq T_3, MV3 = 2) \\ 3 & \text{if } (L3 \geq L3_{\max}, MV3 = 1) \\ 0 & \text{if } (\tau_3 \geq T_3, MV3 = 3) \\ 2 & \text{if } (L3 \leq L3_{\min}, MV3 = 0), \end{cases}$$

for some positive constants  $L3_{\max} \geq L3_{\min} > 0$ . Moreover,  $\tau_3$  resets to 0 each time the PLC switches the value of MV3.

In compact form, the discrete behavior of MV3 and  $\tau_3$  can be modeled by the following constrained difference equation:

$$\begin{pmatrix} \tau_3^+ \\ MV3^+ \end{pmatrix} = \begin{pmatrix} 0 \\ G_{MV3}(MV3) \end{pmatrix} \quad (L3, \tau_3, MV3) \in D_{MV3},$$

where the set  $D_{MV3}$ , captures the update of MV3 and  $\tau_3$

$$\begin{aligned} D_{MV3} &:= D'_{MV3} \cup D''_{MV3}, \\ D'_{MV3} &:= \{(L3, \tau_3, MV3) : L3 \leq L3_{\min}, MV3 = 0\} \cup \\ &\quad \{(L3, \tau_3, MV3) : L3 \geq L3_{\max}, MV3 = 1\}, \\ D''_{MV3} &:= \{(L3, \tau_3, MV3) : \tau_3 \geq T_3, MV3 \in \{2, 3\}\}, \end{aligned}$$

and the function  $G_{MV3}$ , capturing the update law of MV3, is given by

$$G_{MV3}(MV3) := \begin{cases} 3 - MV3 & \text{if } (L3, \tau_3, MV3) \in D''_{MV3} \\ MV3 + 2 & \text{if } (L3, \tau_3, MV3) \in D'_{MV3}. \end{cases}$$

We also found that unless P3 is turned off by the remote SCADA operator, it is always on.

(A1) By default, the system is constantly delivering water; namely, P3 is always equal to 1.

**Stage 2.** Controller C2 actuates the motor valve and the pump in Stage 2 using the control signals MV2 and P2, respectively. Such a decision is based on the information received from a local sensor measuring the water level L2 in Stage 2 and remote information provided by the controller C3 concerning the level of the water L3 in the third stage. The decision rules governing P2 are as follows:

$$P2 := \begin{cases} 1 & \text{if } (L3 \leq L3_{\min}, P2 = 0) \vee \\ & (MV3 \in \{1, 2\}, P2 = 0), \\ 0 & \text{if } (L3 \geq L3_{\max}, P2 = 1) \vee \\ & (MV3 \in \{0, 3\}, P2 = 1). \end{cases} \quad (16)$$

In compact form, the behavior of P2 can be modeled by the following constrained difference equation:

$$P2^+ = G_{P2}(P2) \quad (P2, L3, MV3) \in D_{P2},$$

where  $G_{P2}(P2) := 1 - P2$  and

$$\begin{aligned} D_{P2} &:= \{(P2, L3, MV3) : L3 \leq L3_{\min}, P2 = 0\} \cup \\ &\quad \{(P2, L3, MV3) : MV3 \in \{1, 2\}, P2 = 0\} \cup \\ &\quad \{(P2, L3, MV3) : L3 \geq L3_{\max}, P2 = 1\} \cup \\ &\quad \{(P2, L3, MV3) : MV3 \in \{0, 3\}, P2 = 1\}. \end{aligned}$$

The motor valve in Stage 2 should go through a transition step when migrating from being 1 to 0 and vice versa. Each transition mode has a duration denoted  $T_2 > 0$  seconds. To model this time delay, a timer variable  $\tau_2 \in [0, T_2]$  is used. As a consequence, the decision rules governing the behavior of MV2 are as follows:

$$MV2 := \begin{cases} 2 & \text{if } (L2 \leq L2_{\min}, MV2 = 0) \\ 1 & \text{if } (\tau_2 \geq T_2, MV2 = 2) \\ 3 & \text{if } (L2 \geq L2_{\max}, MV2 = 1) \\ 0 & \text{if } (\tau_2 \geq T_2, MV2 = 3), \end{cases} \quad (17)$$

for some positive constants  $L2_{\max} \geq L2_{\min} > 0$ . Moreover, we switch the value of  $\tau_2$  to 0 (timer reset) each time MV2 switches.

In a compact form, the discrete behavior of MV2 and  $\tau_2$  can be modeled by the following constrained difference equation:

$$\begin{bmatrix} \tau_2^+ \\ MV2^+ \end{bmatrix} = \begin{bmatrix} 0 \\ G_{MV2}(MV2) \end{bmatrix} \quad (L2, \tau_2, MV2) \in D_{MV2},$$

where  $D_{MV2} := D'_{MV2} \cup D''_{MV2}$ ,

$$\begin{aligned} D'_{MV2} &:= \{(L2, \tau_2, MV2) : L2 \geq L2_{\max}, MV2 = 1\} \cup \\ &\quad \{(L2, \tau_2, MV2) : L2 \leq L2_{\min}, MV2 = 0\}, \\ D''_{MV2} &:= \{(L2, \tau_2, MV2) : \tau_2 \geq T_2, MV2 \in \{2, 3\}\}, \end{aligned}$$

and  $G_{MV2}(MV2) := \begin{cases} 3 - MV2 & \text{if } (L2, \tau_2, MV2) \in D''_{MV2} \\ MV2 + 2 & \text{if } (L2, \tau_2, MV2) \in D'_{MV2}. \end{cases}$

**Stage 1.** Similarly to Stage 2, on Stage 1, controller C1 actuates the motor valve and the pump in Stage 1 using the control signals MV1 and P1, respectively. Such a decision is based on the information received from a local sensor measuring the water level L1 in Stage 1 and remote information provided by the controller C2 concerning the level of the water L2 in the second stage. The decision rules governing P1 are as follows:

$$P1 := \begin{cases} 1 & \text{if } (L2 \leq L2_{\min}, P1 = 0) \vee \\ & (MV2 \in \{1, 2\}, P1 = 0), \\ 0 & \text{if } (L2 \geq L2_{\max}, P1 = 1) \vee \\ & (MV2 \in \{0, 3\}, P1 = 1). \end{cases} \quad (18)$$

In compact form, the behavior of P1 can be modeled by the following constrained difference equation:

$$P1^+ = G_{P1}(P1) \quad (P1, L2, MV2) \in D_{P1},$$

where  $G_{P1}(P1) := 1 - P1$  and

$$\begin{aligned} D_{P1} &:= \{(P1, L2, MV2) : L2 \leq L2_{\min}, P1 = 0\} \cup \\ &\quad \{(P1, L2, MV2) : MV2 \in \{1, 2\}, P1 = 0\} \cup \\ &\quad \{(P1, L2, MV2) : L2 \geq L2_{\max}, P1 = 1\} \cup \\ &\quad \{(P1, L2, MV2) : MV2 \in \{0, 3\}, P1 = 1\}. \end{aligned}$$

The motor valve in Stage 1 should go through a transition step when migrating from 1 to 0 and vice versa. Each transition mode has a duration denoted  $T_1 > 0$  seconds. To model this time delay, a timer variable  $\tau_1 \in [0, T_1]$  is used. As a consequence, the decision rules governing the behavior of MV1 are as follows:

$$MV1 := \begin{cases} 2 & \text{if } (L1 \leq L1_{\min}, MV1 = 0) \\ 1 & \text{if } (\tau_1 \geq T_1, MV1 = 2) \\ 3 & \text{if } (L1 \geq L1_{\max}, MV1 = 1) \\ 0 & \text{if } (\tau_1 \geq T_1, MV1 = 3), \end{cases} \quad (19)$$

for some positive constants  $L1_{\max} \geq L1_{\min} > 0$ . Moreover, we switch the value of  $\tau_1$  to 0 (timer reset) each time MV1 switches.

In a compact form, the discrete behavior of MV1 and  $\tau_1$  can be modeled by the following constrained difference equation:

$$\begin{bmatrix} \tau_1^+ \\ MV1^+ \end{bmatrix} = \begin{bmatrix} 0 \\ G_{MV1}(MV1) \end{bmatrix} \quad (L1, \tau_1, MV1) \in D_{MV1},$$

where  $D_{MV1} := D'_{MV1} \cup D''_{MV1}$ ,

$$D'_{MV1} := \{(L1, \tau_1, MV1) : L1 \geq L1_{\max}, MV1 = 1\} \cup \{(L1, \tau_1, MV1) : L1 \leq L1_{\min}, MV1 = 0\},$$

$$D''_{MV1} := \{(L1, \tau_1, MV1) : \tau_1 \geq T_1, MV1 \in \{2, 3\}\},$$

$$\text{and } G_{MV1}(MV1) := \begin{cases} 3 - MV1 & \text{if } (L1, \tau_1, MV1) \in D''_{MV1} \\ MV1 + 2 & \text{if } (L1, \tau_1, MV1) \in D'_{MV1}. \end{cases}$$

While we extracted the discrete model directly, we also developed an automated tool to help us with this analysis. Details are in Appendix A.

## 5.2. Physical Invariants

We now turn our attention to the continuous dynamics. In the absence of attacks, the rate of change of the water level  $L3$  in Stage 3 depends only on the values of MV3 and P3, which coincides with the actual states of the motor valve and the pump in Stage 3:

$$\dot{L3} = F_{L3}(MV3, P3), \quad (20)$$

for some  $F_{L3} : \{0, 1, 2, 3\} \times \{0, 1\} \rightarrow \mathbb{R}$  satisfying the following properties:

$$\begin{aligned} F_{L3}(2, 1) &= F_{L3}(3, 1) > 0, \\ F_{L3}(0, 1) &< 0, \quad F_{L3}(1, 1) > 0. \end{aligned} \quad (21)$$

Similarly, the rate of change of the water level  $L2$  in Stage 2 satisfies

$$\dot{L2} = F_{L2}(MV2, P2), \quad (22)$$

for some  $F_{L2} : \{0, 1, 2, 3\} \times \{0, 1\} \rightarrow \mathbb{R}$  satisfying

$$\begin{aligned} F_{L2}(0, 0) &= 0, \quad F_{L2}(2, 0) = F_{L2}(3, 0) > 0, \\ F_{L2}(0, 1) &< 0, \quad F_{L2}(2, 1) = F_{L2}(3, 1) < 0, \\ F_{L2}(1, P2) &> 0 \quad \forall P2 \in \{0, 1\}. \end{aligned} \quad (23)$$

Furthermore, since we are not considering attacks in this section, the dynamics of  $L2$  do not depend on MV3, which means that when the pump  $P2 = 1$  and the motor valve  $MV3 = 2$  or 1, the same water stream is removed from the second tank due to the length of the channel

between P2 and MV3. Furthermore, according to (16), for each  $i \in \{1, 2\}$ , the pump  $Pi$  and the motor valve  $MVi+1$  are such that

$$\begin{aligned} Pi &= 1 \iff MVi+1 \in \{1, 2\} \\ Pi &= 0 \iff MVi+1 \in \{0, 3\}. \end{aligned} \quad (24)$$

This property implies that the stages are cascaded: the behavior of Stage 1 depends only on the variables of Stages 1 and 2, the behavior of Stage 2 depends only on the variables of Stages 2 and 3, and the behavior of Stage 3 depends only on its own variables. In addition to simplifying the analysis, (24) guarantees a safe operation mode for P1 and P2, otherwise when  $MV2 = 0$  and  $P1 = 1$ , the motor valve MV2 can be damaged due to the pressure that P1 imposes.

## 5.3. Hybrid Model

**Stage 3.** The state vector of Stage 3 is  $x_3 := [L3 \ \tau_3 \ MV3]^T \in X_3$  where  $X_3 := \mathbb{R}_{\geq 0} \times [0, T_3] \times \{0, 1, 2, 3\}$ . Furthermore, the hybrid system  $\mathcal{H}_3$  modeling Stage 3 is given by

$$\mathcal{H}_3 : \begin{cases} \dot{x}_3 = F_3(x_3) & x_3 \in C_3 \\ x_3^+ = G_3(x_3) & x_3 \in D_3, \end{cases} \quad (25)$$

where  $C_3 := \text{cl}(X_3 \setminus D_3)$ ,  $D_3 := D_{MV3}$ ,  $G_3(x_3) := [L3 \ 0 \ G_{MV3}(MV3)]^T$ ,  $F_3(x_3) := [F_{L3}(MV3, 1) \ 1 \ 0]^T$ .

**Stage 2.** The state vector of Stage 2 is  $x_2 := [L2 \ \tau_2 \ MV2 \ P2]^T \in X_2$ , where  $X_2 := \mathbb{R}_{\geq 0} \times [0, T_2] \times \{0, 1, 2, 3\} \times \{0, 1\}$ . Furthermore, the behavior of Stage 2 is influenced by variables  $(L3, MV3)$  of Stage 3. Hence, we introduce the disturbance vector  $u_2 := (L3, MV3) \in \mathcal{U}_2 := [L3_{\min}, L3_{\max} + \delta] \times \{0, 1, 2, 3\}$ , for some  $\delta > 0$  to be specified later. As a result, Stage 3 can be modeled by the disturbed hybrid system  $\mathcal{H}_2$  given by:

$$\mathcal{H}_2 : \begin{cases} \dot{x}_2 = F_2(x_2) & (x_2, u_2) \in C_2 \\ x_2^+ = G_2(x_2, u_2) & (x_2, u_2) \in D_2, \end{cases} \quad (26)$$

where  $C_2 := \text{cl}((X_2 \times \mathcal{U}_2) \setminus D_2)$ ,  $D_2 := (\mathbb{R}_{\geq 0} \times [0, T_2] \times \{0, 1, 2, 3\} \times D_{P2}) \cup (D_{MV2} \times \{0, 1\} \times \mathcal{U}_2)$ ,

$$G_2(x_2, u_2) := [L2 \ G_{22}(x_2, u_2) \ G_{23}(x_2, u_2) \ G_{24}(x_2, u_2)]^T,$$

$$G_{22}(x_2, u_2) := \begin{cases} 0 & \text{if } (L2, \tau_2, MV2) \in D_{MV2} \\ \tau_2 & \text{otherwise,} \end{cases}$$

$$G_{23}(x_2, u_2) := \begin{cases} G_{MV2}(MV2) & \text{if } (L2, \tau_2, MV2) \in D_{MV2} \\ MV2 & \text{otherwise,} \end{cases}$$

$$G_{24}(x_2, u_2) := \begin{cases} G_{P2}(P2) & \text{if } (P2, u_2) \in D_{P2} \\ P2 & \text{otherwise,} \end{cases}$$

$$F_2(x_2) := [F_{L2}(MV2, P2) \ 1 \ 0 \ 0]^T.$$

**Stage 1.** The hybrid equation  $\mathcal{H}_1 := (C_1, F_1, D_1, G_1)$  modeling the first stage is given by:

$$\mathcal{H}_1 : \begin{cases} \dot{x}_1 = F_1(x_1) & (x_1, u_1) \in C_1 \\ x_1^+ = G_1(x_1, u_1) & (x_1, u_1) \in D_1, \end{cases} \quad (27)$$

where  $x_1 := [L1 \ \tau_1 \ MV1 \ P1]^T \in X_1 := \mathbb{R}_{\geq 0} \times [0, T_1] \times \{0, 1, 2, 3\} \times \{0, 1\}$  is the state vector of Stage 1, and  $u_1 := (L2, MV2) \in \mathcal{U}_1 := [L2_{\min} - \delta, L2_{\max} + \delta] \times \{0, 1, 2, 3\}$

is the disturbance vector formed by variables from Stage 1. Finally, the data  $(C_1, F_1, D_1, G_1)$  is constructed the same way as  $(C_2, F_2, D_2, G_2)$ .

**Remark 1.** The behavior of the process is fully represented by the state vector  $x := (x_1, x_2, x_3)$  (and associated numerical values in Appendix B) where  $x_3 := (L3, \tau_3, MV3) \in X_3$ ,  $x_2 := (L2, \tau_2, MV2, P2) \in X_2$ ,  $x_1 := (L1, \tau_1, MV1, P1) \in X_1$ . This is the first accurate and formal model of an industrial system that has been used in the past in various security conferences [63], [12], [21], [9], [3]. We believe this white-box model (in the sense that everything can be explained from first principles, as opposed to a black box model producing outputs from inputs without any explanation of their relationship) can help future researchers extend and improve their security studies for SWaT.

**Remark 2.** In addition to creating the most comprehensive formal description of this popular process, we are also the first to prove that the system is safe without attacks. To analyze the safety of the water treatment plant, we consider a solution  $x := (x_1, x_2, x_3)$  starting from the initial set

$$X_o := X_{o1} \times X_{o2} \times X_{o3}, \quad (28)$$

$$X_{o1} := \{x_1 \in X_1 : L1 \in [L1min, L1max]\}, \quad (29)$$

$$X_{o2} := \{x_2 \in X_2 : L2 \in [L2min, L2max]\}, \quad (30)$$

$$X_{o3} := \{x_3 \in X_3 : L3 \in [L3min, L3max]\}. \quad (31)$$

It is important to note that, due to the transition modes delaying the reaction of the motor valves, it is not possible to guarantee that such a solution  $x$  remains in the set  $X_o$ . However, we will be able to show that such a solution  $x$  remains in a larger set

$$X_s := X_{s1} \times X_{s2} \times X_{s3}, \quad (32)$$

where, for some  $\delta > 0$  to be quantified, we have

$$X_{s1} := \{x_1 \in X_1 : L1 \in [L1min - \delta, L1max + \delta]\}, \quad (33)$$

$$X_{s2} := \{x_2 \in X_2 : L2 \in [L2min - \delta, L2max + \delta]\}, \quad (34)$$

$$X_{s3} := \{x_3 \in X_3 : L3 \in [L3min, L3max + \delta]\}. \quad (35)$$

As a consequence, for any unsafe set  $X_u$  satisfying  $X_u \subset (X_1 \setminus X_{s1}) \times (X_2 \setminus X_{s2}) \times (X_3 \setminus X_{s3})$ , we show that the plant is safe with respect to  $(X_o, X_u)$ .

Our detailed proof of safety in the absence of attacks can be found in Appendix D. From the lengthy proof in the Appendix, it is clear that these mathematical arguments are non-trivial. As far as we are aware, this is the first time this popular water treatment process has been proven safe (without attacks).

## 6. Security Proofs Under Actuation Attacks

Modeling actuator attacks in SWaT is not easy nor simple. Therefore, the details of the full model of SWaT under actuation attacks are presented in Appendix E. Using this attack model, we now analyze the effect of cyber-attacks and study if our new proposed countermeasures can make the system more resilient. In this section, we sketch our results, but the full proofs of adversarial safety can be found in Appendix F.

In the presence of attacks, the variables  $(P_i, MV_i)$  do not necessarily correspond to the actual states of the  $i$ -th motor valve and the  $i$ -th pump, respectively. For this reason, we introduce the extra variables  $(MV1^a, MV2^a, MV3^a)$  to denote the actual states of the motor valves,  $(\tau_1^a, \tau_2^a, \tau_3^a)$  to time the actual transitions of the motor valves, and  $(P1^a, P2^a, P3^a)$  to denote the actual states of the pumps.

As we will show in the next section (experimental results), the system with the original PLC programs is unsafe in the presence of attacks. Indeed, due to the constant demand of water by Stage 3,  $L2$  becomes less than  $L2min - \delta$  if  $P1^a = 0$  is maintained by the attacker. Similarly,  $L3$  becomes less than  $L3min$  if  $P2^a = 0$  is maintained. Therefore the original system is unsafe to attacks that can compromise either the first or the second pump. However, as we will show in this section, if we change the control logic of PLCs, the system can be made safe against arbitrary attacks (as long as they compromise only one control signal).

**Theorem 2.** Consider the hybrid system  $\mathcal{H}_2$ . Consider the initial set  $X_{o2}$  in (30) and the unsafe set  $X_{u2} \subset X_2 \setminus X_{s2}$  with  $X_{s2}$  introduced in (34). Assume that there exist  $\sigma_h > 0$  and  $\sigma_g > 0$  such that

$$4T_2(F_{L2}(2, 0, 0) + \sigma_h) \leq \delta, \quad (36)$$

$$4T_2(F_{L2}(2, 1, 1) + \sigma_g) \leq \delta. \quad (37)$$

Then, the hybrid system  $\mathcal{H}_2$  is safe with respect to  $(X_{o2}, X_{u2})$  uniformly in  $(u_2, w_2) \in \mathcal{U}_2 \times \mathcal{W}_2$ , and admits a barrier function certificate given by

$$B(x_2) := (L2 - L2min + g(\tau_2, MV2)) \times (L2 - L2max - h(\tau_2, MV2)),$$

where  $g(\tau_2, MV2) := (-F_{L2}(2, 1, 1) + \sigma_g) * [\tau_2 + T_2 * w_g(MV2)]$ ,  $w_g(3) := 0$ ,  $w_g(0) := 1$ ,  $w_g(2) := 2$ ,  $w_g(1) := 3$ ,  $h(\tau_2, MV2) := (F_{L2}(2, 0, 0) + \sigma_h) * [\tau_2 + T_2 * w_h(MV2)]$ ,  $w_h(2) := 0$ ,  $w_h(1) := 1$ ,  $w_h(3) := 2$ ,  $w_h(0) := 3$ .  $\square$

We now harden the system to make it more resilient to attacks. In particular, we first change the control logic of the PLC controlling stage 3 ( $C3$ ) so that  $P3$  is not always 1. As a result, we include  $P3$  as a control parameter governed by the following logic:

$$P3 := \begin{cases} 0 & \text{if } (L3 \leq L3o, P3 = 1) \\ 1 & \text{if } (L3 \geq L3o, P3 = 0), \end{cases} \quad (38)$$

where  $L3o > 0$  is a lower bound on the water level  $L3$  in Stage 3, it aims to avoid the dry-runs (operates without

liquid) of the pump P3. Hence, the behavior of P3 can be modeled by the following constrained difference equation:

$$P3^+ = G_{P3}(P3) \quad (L3, P3) \in D_{P3},$$

where  $D_{P3} := \{(L3, P3) : L3 \leq L3_o, P3 = 1\} \cup \{(L3, P3) : L3 \geq L3_o, P3 = 0\}$  and  $G_{P3}(P3) := 1 - P3$ .

Using the logic (84) with  $L3_o = L3_{min}$ , we are able to show the following claim.

**Claim 1.** When P3 is governed by (84) with  $L3_o = L3_{min}$ , the plant remains safe under any arbitrary time series of possible attacks affecting  $P2^a$ . •

To show Claim 5, we use Theorem 11 to conclude that it is enough to show the safety of Stage 3 uniformly in  $P2^a \in \{0, 1\}$  when P3 is governed by (84). To simplify the analysis, we model Stage 3 when only  $P2^a$  is attacked and P3 is controlled by our modified control logic (it is hardened); see system  $\mathcal{H}_3$ .

**Theorem 3.** Consider the hybrid system  $\mathcal{H}_3$ . Consider the initial set  $\bar{X}_{o3} := X_{o3} \times \{0, 1\}$  and an unsafe set  $X_{u3} \subset \bar{X}_3 \setminus \bar{X}_{s3}$  with  $\bar{X}_{s3} = X_{s3} \times \{0, 1\}$ . Assume that there exists  $\sigma > 0$  such that

$$4T_3(F_{L3}(3, 0, 1) + \sigma) \leq \delta. \quad (39)$$

Then, the hybrid system  $\mathcal{H}_3$  is safe with respect to  $(\bar{X}_{o3}, \bar{X}_{u3})$  uniformly in  $u_3 = w_{2m} = P2^a \in \mathcal{U}_3$ , and admits a barrier function certificate given by

$$B(\bar{x}_3) := (L3 - L3_{min})(L3 - L3_{max} - P3 * f(\tau_3, MV3)),$$

where  $f(\tau_3, MV3) := (F_{L3}(3, 0, 1) + \sigma)[\tau_3 + T_3 * w_f(MV3)]$ , and  $w_f(2) := 0$ ,  $w_f(1) := 1$ ,  $w_f(3) := 2$ ,  $w_f(0) := 3$ . □

So far, we showed that our control logic modification makes the system safe against attacks in P2, but in the next section on experimental results, we show that our change is not enough when the adversary attacks P1. As a result, we need to modify the PLC controlling stage 2 as well, i.e., C2.

**Claim 2.** When modifying the logic in (16) governing P2 as follows:

$$P2 := \begin{cases} 1 & \text{if } (L3 \leq L3_{min}, P2 = 0, L2 \geq L2_{min}) \vee \\ & (MV3 \in \{1, 2\}, P2 = 0, L2 \geq L2_{min}), \\ 0 & \text{if } (L3 \geq L3_{max}, P2 = 1) \vee \\ & (MV3 \in \{0, 3\}, P2 = 1) \vee \\ & (L2 \leq L2_{min}, P2 = 1), \end{cases} \quad (40)$$

the plant becomes safe under any attack affecting  $P1^a$ . •

Since only  $P1^a$  is attacked, the model of Stage 3 is as in (25), and its safety is already analyzed in Theorem 6. Hence, to prove Claim 6, it is enough to prove that Stage 2 is safe uniformly in  $(P1^a, x_3) \in \{0, 1\} \times X_{s3}$  when only  $P1^a$  is attacked and when (87) governs P2. To simplify the proof, we model Stage 2 when only  $P1^a$  is attacked and when the logic governing P2 is modified;

$$\mathcal{H}_2 : \begin{cases} \dot{x}_2 = F_2(x_2, u_2) & (x_2, u_2) \in \tilde{C}_2 \times \{0, 1\} \\ x_2^+ = \tilde{G}_2(x_2, L3, MV3) & (x_2, u_2) \in \tilde{D}_2 \times \{0, 1\}, \end{cases} \quad (41)$$

where

$$u_2 := (L3, MV3, P1^a) \in \mathcal{U}_2,$$

$$\mathcal{U}_2 := [L3_{min}, L3_{max} + \delta] \times \{0, 1, 2, 3\} \times \{0, 1\},$$

and

$$F_2(x_2, u_2) := \begin{bmatrix} F_{L2}(MV2, P2, P1^a) \\ 1 \\ 0 \\ 0 \end{bmatrix},$$

**Theorem 4.** Consider the hybrid system  $\mathcal{H}_2$  in (88). Consider the initial set  $X_{o2}$  in (30) and the unsafe set  $X_{u2} \subset X_2 \setminus X_{s2}$  with  $X_{s2}$  introduced in (34). Assume that there exists  $\sigma_h > 0$  such that

$$4T_2(F_{L2}(2, 0, 1) + \sigma_h) \leq \delta. \quad (42)$$

Then, the hybrid system  $\mathcal{H}_2$  in (88) is safe with respect to  $(X_{o2}, X_{u2})$  uniformly in  $u_2 \in \mathcal{U}_2$ , and admits a barrier function certificate given by

$$B(x_2) := (L2 - L2_{min})(L2 - L2_{max} - \chi(L2) * h(\tau_2, MV2)), \quad (43)$$

where  $\chi : \mathbb{R} \rightarrow [0, 1]$  is a smooth function such that

$$\begin{cases} \chi(L2) = 1 & \text{if } L2 \geq L2_{max} \\ \chi(L2) = 0 & \text{if } L2 \leq L2_{min} \\ \chi(L2) \in [0, 1] & \text{otherwise,} \end{cases}$$

and  $h(\tau_2, MV2) := (F_{L2}(2, 0, 1) + \sigma_h) * [\tau_2 + T_2 * w_h(MV2)]$ ,  $w_h(2) := 0$ ,  $w_h(1) := 1$ ,  $w_h(3) := 2$ ,  $w_h(0) := 3$ . □

In summary, we have shown that the original SWaT system cannot guarantee safety when the attacker compromises any of the following actuators P1, P2, MV1, MV2, MV3. However, we proposed a set of control logic changes to PLCs 1 and 2, and with these changes, we were able to prove that the system is safe if the attacker compromises any of these actuators: P1, P2, P3, MV2, MV3. The only time SWaT cannot guarantee safety is when the attacker compromises MV1 or compromises more than one actuator. The reason for this is that the amount of water coming into the first tank is higher than the amount that can be taken out by P1. To guarantee safety against a compromise of MV1 we would need a physical redesign of the system so that the rate of flow of entering water is the same as the rate of flow that P1 can take out of the first tank.

## 7. Experimental Results

Our experiments in the real-world system confirm the theoretical results in Section 6. As depicted in Figs. 6 and 7, the original SWaT system reaches an unsafe state (L below 750mm) showing the attacks  $P1^a = 0$  and  $P2^a = 0$  were effective.

As stated in Section F.1, PLC programs can be modified to make the system more resilient against actuator attacks. We update the PLC program, including additional

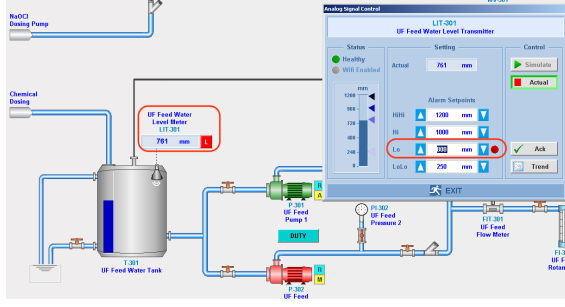


Figure 6: SCADA reports that the tank level is out of its operational limits; see the red box with the **L** character. In the detailed view, a red circle shows which boundary is being transgressed.

validations to the control logic. The program slice in Listing 1 shows how we code Equations (87) into the PLC program. To know how to change the PLC program, we need to understand the program semantics at a high level. The PLC program uses a bit array variable ( $P2.SD$  in Listing 1) below to guarantee that the actuators always operate within the safety conditions; the code computes the safety conditions between lines 7 and 15, each one assigns to a different position in the array. Later in the code, the **IF** statement at line 17 checks that all safety conditions hold before turning ON the pump. To code the additional conditions, we translate Equation (87) into logical expression and insert them between lines 7 and 15.

```

1  VAR
2      P2.SD           : ARRAY[0..15] OF BOOL;
3      P2.Auto, P2.Fault, P2.Permissive : BOOL;
4      P2.FT_Start, P2.FT_Stop, P2.Start : BOOL;
5  END VAR
6  ...
7  P2.SD.0 := LS2.Alarm;
8  P2.SD.1 := P2.STATUS=2 AND MV2.STATUS<>2;
9  P2.SD.2 := F2_TM.DN;
10 P2.SD.3 := L1.Level < L1min;
11 P2.SD.4 := 0;
12 ...
13 P2.SD.15:= 0;
14 ...
15 IF P2.Auto THEN
16     IF NOT P2.Fault AND NOT P2.FT_Start
17       AND NOT P2.FT_Stop AND (P2.Permissive==1)
18       AND P2.SD=0 THEN
19         P2.Start :=1;           (*Turn ON P2*)
20     ELSE IF P2.Start OR (P2.SD<>0) OR P2.Fault THEN
21         P2.Start :=0;           (*Turn OFF P2*)
22     END IF;
23 END IF;
24 ...

```

Listing 1: Slice of hardened PLC program controlling P2.

We test the enhanced versions of controllers with the whole set of pump attacks. Fig. 7 (bottom) shows how the controllers mitigate the effects of the pump attacks. While the same attacks led the system to unsafe states under unprotected controllers, the enhanced version of the PLC programs allows the system to respond to the attacks effectively. Again, this matched our theoretical results in the previous section.

We now turn our attention to the fidelity of our model to the real-world system. We compare our model to traces

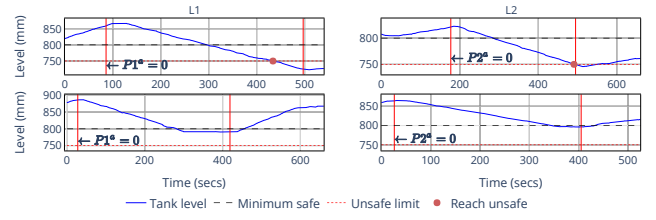


Figure 7: Response to pump attacks ( $P1^a$  and  $P2^a$ ). Solid red lines show when the attack starts and when it ends. Top: Tank levels reach unsafe states. Bottom: With enhanced PLC programs, controllers change the strategy when the tanks are under the lower limit.

from the real-world operation of SWaT and also to a previously proposed simulation of SWaT from a paper from IEEE S&P 2018 [12]. Fig. 8 shows how our proposed model closely follows the real-world operation of the system, while the previous simulation differs significantly from the real-world operation. In the Appendix (see Fig. 13), we include additional simulations of the proposed model to evaluate the correctness of the model empirically. We also emphasize that our implementation is based on the equations provided in this paper, while the previous simulation does not have equations for the behavior of the system, so this previous work cannot be used to reason mathematically about the safety properties of the system.

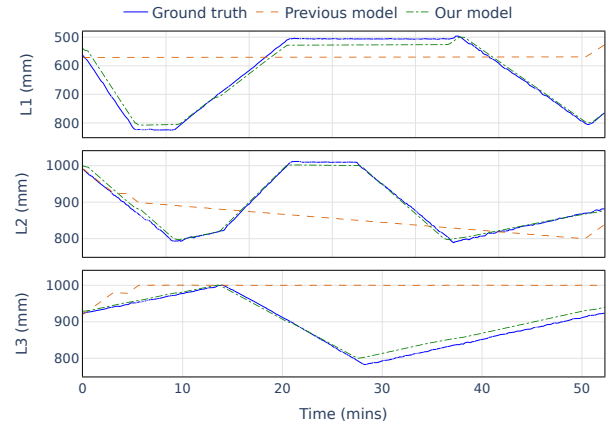


Figure 8: Comparing our SWaT model with the SWaT simulator [12] (previous model) for a period of 50 minutes.

We finalize this experimental section by discussing some practicalities we discovered while trying to launch the attacks in SWaT. In our scenario, we have two pairs of Pump motor-valve in-line connections. In this setup, if the pump opens while the motor valve is closed, it will cause a dead-head effect. To mitigate the undesired effect, controllers check the motor valve state before opening the pumps. When we launch the attack (open the pump) using the SCADA, the controller blocks the action because the next motor valve is closed. To get around this challenge, we modify the PLC code to overwrite the validation and force the output as desired. A realistic attacker can produce the same effect via a Man-in-the-middle attack, spoofing the data from the controller to the pump.

Another practical challenge we found is that after launching our attacks, the SCADA triggered an alarm linked to the pump (see Fig. 9). This is due to inconsistency. While MV is closed, the controller tries to close P, but the attacker swaps the action to open. When the SCADA reads the pump state, it reports an open state, causing the discrepancy. The attacker can hide this effect by spoofing the actuator state read back to the SCADA.

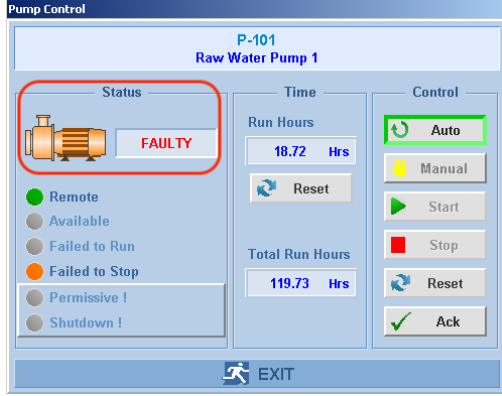


Figure 9: Effect of the attack  $P1^a = 1$ . The PLC detects pump is not OFF after sending OFF the command, so it sends this alert to the SCADA server

## 8. Discussion

**Limitations:** One limitation of our approach is that the effort in modeling real-world systems and then proving safety can be significant. Section 4 shows that analyzing two simple examples (a thermostat and a robotic vehicle) is straightforward; however, analyzing SWaT requires considerable effort. Our lengthy Appendix is an indication of the several months it took us to (1) model SWaT with hybrid mathematical equations, (2) model SWaT under actuation attacks, and (3) prove (or disprove with a counterexample) that the system is safe.

Unfortunately, there is no free lunch. Proving safety in cyber-physical systems is generally a very complex process, and state-of-the-art tools are not scalable to complex systems with hundreds or thousands of actuators (e.g., a power grid). However, several real-world control systems (like SWaT) are in the range ( $\approx 20 - 50$  actuators and sensors) where first-principles models can be created by experts, so while our framework might not be scalable to systems with hundreds of actuators and sensors, we can still model a large amount of practical real-world control systems. Essentially, if a real-world system is similar in scale to SWaT, we can use our formal verification approaches.

**Generalization:** Our framework can be applied to various systems. In this paper, we have used our models to show the insecurity of a robotic vehicle and showed how to change the control algorithms to make this robotic vehicle secure against arbitrary false data injection attacks when the attacker compromises only one control signal. We then showed how to do the same analysis for a real-world water system. In general, we can use our methods in other systems where the use of barrier functions is tractable. These include bipedal robots [44], autonomous robotic systems [25], and drones [56].

**Finding barrier functions:** One of the challenges of our approach is finding barrier functions. Depending on how complex the CPS is and depending on how much understanding we have of the system's behavior, the search for barrier functions can be found intuitively by the engineer by exploiting the physics or decomposing the system into interconnected subsystems. If the system is too complex, the search for barrier functions can be performed numerically. Here, different algorithms such as SOS [48], [66] or learning-based methods [49], [50]. No matter how you look for a barrier function, your conclusions are mathematically rigorous once you find it.

## 9. Conclusion

In this paper, we have presented the most comprehensive formal model of a popular real-world system used by the CPS security community. We have proved for the first time that the system is safe under normal operations and also showed that the system is not safe under actuation attacks. We then showed how to modify the control logic of the PLCs so that no single actuation attack can make the system unsafe (except for MV1). This latter guarantee is proof of security against any attack tactic on any single actuator. Finally, we have shown how compromising more than one actuator results in an unsafe system.

Our main contribution is to push state-of-the-art provable security guarantees in cyber-physical systems. We argued that the progress of formal security guarantees in the past forty years has first needed precise definitions and claims. Precise definitions of security are the building block to proposing refutable assertions, allowing us to follow the scientific method because future papers can contradict or build upon them. So far, the literature on industrial control in security conferences lacks these assertions. Our paper has shown how to use the concept of barrier function certificates to develop these security assertions in a real-world process.

We hope that this work and the detailed models of this paper can help future researchers use SWaT and researchers attempting to verify CPS safety under various attacks.

## Acknowledgments

Research was partially sponsored by NSF CNS-1929410, 1931573 and by the Army Research Office and was accomplished under Grant Number W911NF-20-1-0253. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Army Research Office or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation herein. Research by R. G. Sanfelice partially supported by NSF Grants no. CNS-2039054 and CNS-2111688, by AFOSR Grants no. FA9550-19-1-0169, FA9550-20-1-0238, and FA9550-23-1-0145, by AFRL Grant nos. FA8651-22-1-0017 and FA8651-23-1-0004.

## References

- [1] Sridhar Adepu and Aditya Mathur. Distributed detection of single-stage multipoint cyber attacks in a water treatment plant. In *Proceedings of the 11th ACM on Asia Conference on Computer and Communications Security*, pages 449–460, 2016.
- [2] Chuadhry Mujeeb Ahmed, Martin Ochoa, Jianying Zhou, Aditya P Mathur, Rizwan Qadeer, Carlos Murguia, and Justin Ruths. Noiseprint: Attack detection using sensor and process noise fingerprint in cyber physical systems. In *Proceedings of the 2018 on Asia Conference on Computer and Communications Security*, pages 483–497, 2018.
- [3] Chuadhry Mujeeb Ahmed, Jianying Zhou, and Aditya P Mathur. Noise matters: Using sensor and process noise fingerprint to detect stealthy cyber attacks and authenticate sensors in cps. In *Proceedings of the 34th Annual Computer Security Applications Conference*, pages 566–581, 2018.
- [4] Rajeev Alur, Thao Dang, and Franjo Ivančić. Counter-example guided predicate abstraction of hybrid systems. In *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, pages 208–223. Springer, 2003.
- [5] Aaron D Ames, Samuel Coogan, Magnus Egerstedt, Gennaro Notomista, Koushil Sreenath, and Paulo Tabuada. Control barrier functions: Theory and applications. In *2019 18th European control conference (ECC)*, pages 3420–3431. IEEE, 2019.
- [6] AP. Revenge hacker: 34 months, must repay georgia-pacific \$1m. <https://www.usnews.com/news/louisiana/articles/2017-02-16/revenge-hacker-34-months-must-repay-georgia-pacific-1m>, February 2017.
- [7] J. P. Aubin. *Viability Theory*. Birkhauser Boston Inc., Cambridge, MA, USA, 1991.
- [8] Yulong Cao, Chaowei Xiao, Benjamin Cyr, Yimeng Zhou, Won Park, Sara Rampazzi, Qi Alfred Chen, Kevin Fu, and Zhuoqing Morley Mao. Adversarial Sensor Attack on LiDAR-based Perception in Autonomous Driving. In *Conference on Computer and Communications Security (CCS)*, 2019.
- [9] John H Castellanos, Martín Ochoa, and Jianying Zhou. Finding dependencies between cyber-physical domains for security testing of industrial control systems. In *Proceedings of the 34th Annual Computer Security Applications Conference*, pages 582–594, 2018.
- [10] John Henry Castellanos and Jianying Zhou. A modular hybrid learning approach for black-box security testing of cps. In *International Conference on Applied Cryptography and Network Security*, pages 196–216. Springer, 2019.
- [11] Xin Chen, Erika Abraham, and Sriram Sankaranarayanan. Flow\*: An analyzer for non-linear hybrid systems. In *Computer Aided Verification: 25th International Conference, CAV 2013, Saint Petersburg, Russia, July 13-19, 2013. Proceedings 25*, pages 258–263. Springer, 2013.
- [12] Yuqi Chen, Christopher M Poskitt, and Jun Sun. Learning from mutants: Using code mutation to learn and monitor invariants of a cyber-physical system. In *2018 IEEE Symposium on Security and Privacy (SP)*, pages 648–660. IEEE, 2018.
- [13] Anton Cherepanov. Win32/industry, a new threat for industrial control systems. *White Paper. ESET*, 2017.
- [14] Hongjun Choi, Wen-Chuan Lee, Yousra Aafer, Fan Fei, Zhan Tu, Xiangyu Zhang, Dongyan Xu, and Xinyan Deng. Detecting attacks against robotic vehicles: A control invariant approach. In *Proceedings of the 2018 ACM SIGSAC Conference on Computer and Communications Security, CCS '18*, pages 801–816, New York, NY, USA, 2018. ACM.
- [15] Mauro Conti, Denis Donadel, and Federico Turrin. A survey on industrial control system testbeds and datasets for security research. *IEEE Communications Surveys & Tutorials*, 23(4):2248–2294, 2021.
- [16] Gökçen Yılmaz Dayanikli, Rees R Hatch, Ryan M Gerdes, Hongjie Wang, and Regan Zane. Electromagnetic sensor and actuator attacks on power converters for electric vehicles. In *2020 IEEE Security and Privacy Workshops (SPW)*, pages 98–103. IEEE, 2020.
- [17] Danny Dolev and Andrew Yao. On the security of public key protocols. *IEEE Transactions on information theory*, 29(2):198–208, 1983.
- [18] Laurent Doyen, Goran Frehse, George J Pappas, and André Platzer. Verification of hybrid systems. In *Handbook of Model Checking*, pages 1047–1110. Springer, 2018.
- [19] Raj Gautam Dutta, Xiaolong Guo, Teng Zhang, Kevin Kwiat, Charles Kamhoua, Laurent Njilla, and Yier Jin. Estimation of safe sensor measurements of autonomous system under attack. In *Proceedings of the 54th Annual Design Automation Conference 2017*, page 46. ACM, 2017.
- [20] Raj Gautam Dutta, Feng Yu, Teng Zhang, Yaodan Hu, and Yier Jin. Security for safety: a path toward building trusted autonomous vehicles. In *Proceedings of the International Conference on Computer-Aided Design*, page 92. ACM, 2018.
- [21] Cheng Feng, Venkata Reddy Palleti, Aditya Mathur, and Deepthi Chana. A systematic framework to generate invariants for anomaly detection in industrial control systems. In *2019 Network and Distributed System Security Symposium (NDSS)*.
- [22] Xuan Feng, Qiang Li, Haining Wang, and Limin Sun. Acquisitional rule-based engine for discovering internet-of-things devices. In *27th {USENIX} Security Symposium ({USENIX} Security 18)*, pages 327–341, 2018.
- [23] Ilias Giechaskiel and Kasper Rasmussen. Taxonomy and challenges of out-of-band signal injection attacks and defenses. *IEEE Communications Surveys & Tutorials*, 2019.
- [24] Jairo Giraldo, Sahand Hadizadeh Kafash, Justin Ruths, and Alvaro A Cardenas. Daria: Designing actuators to resist arbitrary attacks against cyber-physical systems. In *2020 IEEE European Symposium on Security and Privacy (EuroS&P)*, pages 339–353. IEEE Computer Society, 2020.
- [25] Paul Glotfelter, Ian Buckley, and Magnus Egerstedt. Hybrid nonsmooth barrier functions with applications to provably safe and composable collision avoidance for robotic systems. *IEEE Robotics and Automation Letters*, 4(2):1303–1310, 2019.
- [26] R. Goebel, R. G. Sanfelice, and A. R. Teel. *Hybrid Dynamical Systems: Modeling, Stability, and Robustness*. Princeton University Press, New Jersey, 2012.
- [27] Oded Goldreich. *Foundations of cryptography: volume 2, basic applications*. Cambridge university press, 2009.
- [28] Ziyang Guo, Dawei Shi, Karl Henrik Johansson, and Ling Shi. Optimal linear cyber-attack on remote state estimation. *IEEE Transactions on Control of Network Systems*, 4(1):4–13, 2017.
- [29] Dina Hadžiosmanović, Robin Sommer, Emmanuele Zamboni, and Pieter H Hartel. Through the eye of the plc: semantic security monitoring for industrial processes. In *Proceedings of the 30th Annual Computer Security Applications Conference*, pages 126–135. ACM, 2014.
- [30] Changhua He, Mukund Sundararajan, Anupam Datta, Ante Derek, and John C Mitchell. A modular correctness proof of ieee 802.11 i and tls. In *Proceedings of the 12th ACM conference on Computer and communications security*, pages 2–15. ACM, 2005.
- [31] Thomas A Henzinger. Hybrid automata with finite bisimulations. In *International Colloquium on Automata, Languages, and Programming*, pages 324–335. Springer, 1995.
- [32] Andreas Hoehn and Ping Zhang. Detection of replay attacks in cyber-physical systems. In *2016 American Control Conference (ACC)*, pages 290–295. IEEE, 2016.
- [33] Pushpak Jagtap, Sadegh Soudjani, and Majid Zamani. Temporal logic verification of stochastic systems using barrier certificates. In *International Symposium on Automated Technology for Verification and Analysis*, pages 177–193. Springer, 2018.
- [34] Karl-Heinz John and Michael Tiegelkamp. *IEC 61131-3: programming industrial automation systems: concepts and programming languages, requirements for programming systems, decision-making aids*. Springer Science & Business Media, 2010.
- [35] Blake Johnson, Dan Caban, Marina Krotofil, Dan Scali, Nathan Brubaker, and Christopher Glyer. Attackers deploy new ICS Attack Framework” TRITON” and cause operational disruption to critical infrastructure. *Threat Research Blog*, 2017.



- [36] Sahand Hadizadeh Kafash, Jairo Giraldo, Carlos Murguia, Alvaro A Cardenas, and Justin Ruths. Constraining attacker capabilities through actuator saturation. In *2018 Annual American Control Conference (ACC)*, pages 986–991. IEEE, 2018.
- [37] Jonathan Katz and Yehuda Lindell. *Introduction to modern cryptography*. CRC press, 2014.
- [38] Pingfan Kong, Yi Li, Xiaohong Chen, Jun Sun, Meng Sun, and Jingyi Wang. Towards concolic testing for hybrid systems. In *International Symposium on Formal Methods*, pages 460–478. Springer, 2016.
- [39] Robert M Lee, Michael J Assante, and Tim Conway. German steel mill cyber attack. *Industrial Control Systems*, 30:62, 2014.
- [40] Qin Lin, Sridha Adepu, Sicco Verwer, and Aditya Mathur. Tabor: A graphical model-based approach for anomaly detection in industrial control systems. In *Proceedings of the 2018 on Asia Conference on Computer and Communications Security*, pages 525–536, 2018.
- [41] M. Maghenem and R. G. Sanfelice. Characterizations of safety in hybrid inclusions via barrier functions. In *Proceedings of the 22Nd ACM International Conference on Hybrid Systems: Computation and Control*, HSCC '19, pages 109–118, NY, USA, 2019. ACM.
- [42] Mohamed Maghenem and Ricardo G Sanfelice. Sufficient conditions for forward invariance and contractivity in hybrid inclusions using barrier functions. *Automatica*, 124, 2021.
- [43] Yilin Mo and Bruno Sinopoli. Secure control against replay attacks. In *2009 47th Annual Allerton Conference on Communication, Control, and Computing (Allerton)*, pages 911–918. IEEE, 2009.
- [44] Quan Nguyen and Koushil Sreenath. Safety-critical control for dynamical bipedal walking with precise footstep placement. *IFAC-PapersOnLine*, 48(27):147–154, 2015.
- [45] Jaclyn Peiser. A hacker broke into a florida towns water supply and tried to poison it with lye, police said. <https://www.washingtonpost.com/nation/2021/02/09/oldsmar-water-supply-hack-florida/>, February 2021.
- [46] A. Platzer. *Logical analysis of hybrid systems: proving theorems for complex dynamics*. Springer Science & Business Media, 2010.
- [47] André Platzer. The structure of differential invariants and differential cut elimination. *Logical Methods in Computer Science*, 8, 2012.
- [48] S. Prajna, A. Jadbabaie, and G. J. Pappas. A framework for worst-case and stochastic safety verification using barrier certificates. *IEEE Transactions on Automatic Control*, 52(8):1415–1428, 2007.
- [49] Alexander Robey, Haimin Hu, Lars Lindemann, Hanwen Zhang, Dimos V Dimarogonas, Stephen Tu, and Nikolai Matni. Learning control barrier functions from expert demonstrations. In *2020 59th IEEE Conference on Decision and Control (CDC)*, pages 3717–3724. IEEE, 2020.
- [50] Alexander Robey, Lars Lindemann, Stephen Tu, and Nikolai Matni. Learning robust hybrid control barrier functions for uncertain systems. *IFAC-PapersOnLine*, 54(5):1–6, 2021.
- [51] Imran Sajjad, Daniel D Dunn, Rajnikant Sharma, and Ryan Gerdes. Attack mitigation in adversarial platooning using detection-based sliding mode control. In *Proceedings of the First ACM Workshop on Cyber-Physical Systems-Security and/or Privacy*, pages 43–53, 2015.
- [52] Ricardo Sanfelice, David Copp, and Pablo Nanez. A toolbox for simulation of hybrid systems in matlab/simulink: Hybrid equations (hyeq) toolbox. In *Proceedings of the 16th international conference on Hybrid systems: computation and control*, pages 101–106, 2013.
- [53] Ricardo G Sanfelice. Analysis and design of cyber-physical systems: a hybrid control systems approach. *Cyber-Physical Systems*, pages 3–31, 2015.
- [54] Neetesh Saxena, Leilei Xiong, Victor Chukwuka, and Santiago Grijalva. Impact evaluation of malicious control commands in cyber-physical smart grids. *IEEE Transactions on Sustainable Computing*, 2018.
- [55] Jayaprakash Selvaraj, Gökçen Y Dayanıklı, Neelam Prabhu Gaunkar, David Ware, Ryan M Gerdes, Mani Mina, et al. Electromagnetic induction attacks against embedded systems. In *Asia Conference on Computer and Communications Security (AsiaCCS)*, pages 499–510. ACM, 2018.
- [56] Andrew Singletary, Aiden Swann, Yuxiao Chen, and Aaron Ames. Onboard safety guarantees for racing drones: High-speed geofencing with control barrier functions. *IEEE Robotics and Automation Letters*, 2022.
- [57] Jill Slay and Michael Miller. Lessons learned from the maroochy water breach. In *Critical Infrastructure Protection*, volume 253/2007, pages 73–82. Springer Boston, November 2007.
- [58] Yun Mok Son, Ho Cheol Shin, Dong Kwan Kim, Young Seok Park, Ju Hwan Noh, Ki Bum Choi, Jung Woo Choi, and Yong Dae Kim. Rocking drones with intentional sound noise on gyroscopic sensors. In *USENIX Security Symposium (USENIX Security)*. USENIX Association, 2015.
- [59] Mingshun Sun, Ali Al-Hashimi, Ming Li, and Ryan Gerdes. Impacts of constrained sensing and communication based attacks on vehicular platoons. *IEEE Transactions on Vehicular Technology*, 69(5):4773–4787, 2020.
- [60] Rui Tan, Varun Badrinath Krishna, David KY Yau, and Zbigniew Kalbarczyk. Impact of integrity attacks on real-time pricing in smart grids. In *Proceedings of the 2013 ACM SIGSAC conference on Computer & communications security*, pages 439–450. ACM, 2013.
- [61] Richard J Thomas and Tom Chothia. Learning from vulnerabilities-categorising, understanding and detecting weaknesses in industrial control systems. In *Computer Security: ESORICS 2020 International Workshops, CyberICPS, SECPRE, and ADIoT, Guildford, UK, September 14–18, 2020, Revised Selected Papers 6*, pages 100–116. Springer, 2020.
- [62] Timothy Trippel, Ofir Weisse, Wenyuan Xu, Peter Honeyman, and Kevin Fu. Walnut: Waging doubt on the integrity of mems accelerometers with acoustic injection attacks. In *European Symposium on Security and Privacy (EuroS&P)*, pages 3–18. IEEE, 2017.
- [63] David I Urbina, Jairo A Giraldo, Alvaro A Cardenas, Nils Ole Tippenhauer, Junia Valente, Mustafa Faisal, Justin Ruths, Richard Candell, and Henrik Sandberg. Limiting the impact of stealthy attacks on industrial control systems. In *Conference on Computer and Communications Security (CCS)*, pages 1092–1105. ACM, 2016.
- [64] David I Urbina, Jairo Alonso Giraldo, Nils Ole Tippenhauer, and Alvaro A Cardenas. Attacking fieldbus communications in ics: Applications to the swat testbed. In *SG-CRC*, pages 75–89, 2016.
- [65] Mathy Vanhoef and Frank Piessens. Key reinstallation attacks: Forcing nonce reuse in wpa2. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, pages 1313–1328. ACM, 2017.
- [66] Li Wang, Dongkun Han, and Magnus Egerstedt. Permissive barrier certificates for safe stabilization using sum-of-squares. In *2018 Annual American Control Conference (ACC)*, pages 585–590. IEEE, 2018.
- [67] Yong Wang, Zhaoyan Xu, Jialong Zhang, Lei Xu, Haopei Wang, and Guofei Gu. Srid: State relation based intrusion detection for false data injection attacks in scada. In *European Symposium on Research in Computer Security*, pages 401–418. Springer, 2014.
- [68] C Yan, H Shin, C Bolton, W Xu, Y Kim, and K Fu. Sok: A minimalist approach to formalizing analog sensor security. In *2020 IEEE Symposium on Security and Privacy (SP)*, pages 480–495.
- [69] Kim Zetter. *Countdown to Zero Day: Stuxnet and the launch of the world's first digital weapon*. Broadway books, 2014.
- [70] Lin Zhang, Xin Chen, Fanxin Kong, and Alvaro A Cardenas. Real-time attack-recovery for cyber-physical systems using linear approximations. In *2020 IEEE Real-Time Systems Symposium (RTSS)*, pages 205–217. IEEE, 2020.

## A. Control Invariants from PLC Code

Programmable Logic Controllers, also known as PLCs, are industrial computers with multiple hardware modules to measure and control physical systems. PLCs rely on robust hardware that allows them to operate under

extreme conditions of temperature, moisture, etc. Engineers program multiple routines that mandate how the PLC response to different states of the system. The IEC-61131-3 standard [34] groups four different languages that are broadly used by different vendors to code the routines, the most popular are Ladder Logic and Structured Text. The group of programs and routines that run on the PLCs are denoted as *control logic*.

---

**Algorithm 1:** Extract control conditions

---

**Data:** A PLC source code SC

**Input:** A set of actuators A

**Result:** Returns a dictionary with control conditions for each actuator state.

---

```

1 CFG:= buildCFG(SC) ;      // Build CFG
  from PLC source code
2 for a ∈ A do
3   d:= SC.GetDefinitions(a);
4   D.push(d);
5 R:= hashmap() ;           // Condition set
6 for d ∈ D do
7   b:= CFG.GetBlockIndex(d) ; // Get
  block ID for definition d
8   p:= CFG.Path(b);
9   s:= SC.symbExec(p) ;      // Get
  symbolic path condition for
  path p
10  R.update(d) := s;
11 return R

```

---

PLCs operate reliably following a cyclic pattern called scan cycle. Roughly speaking, (1) stores data from sensor modules to a local buffer, (2) updates the network modules with local buffers, (3) runs the control logic, (4) updates local output buffers to actuator modules, (5) executes internal safety checks, and (6) repeats the cycle.

Consider the slice of PLC code shown in Listing 2. It encodes the control logic to manage a motor valve called MV2. We split the code into basic blocks **B1**, **B2**, ..., **B16** of continuous statements to explain the logic behind the PLC program.

The program starts collecting data from sensors and saving into local variables (B1). PLC programs use special functions like latches, counters and timers, for example in B3, the code employs a latch (MV2.SR) to collect the status of the L2 level, the function SETD evaluates the inputs of the latch and updates the output according to, if the input MV2.SR.S is True the output MV2.SR.Out sets True, otherwise, if the input MV2.SR.R is True the output sets False. Timers use the TONR function. When the code calls TONR, the PLC evaluates if the timer is enabled (Enable input), if so, an internal timer is set to the PRE input value (in milliseconds), this timer is independent of the scan cycle to enforce real-time responses. The program has two timers (B4), one for each transition of MV2 OFF to ON and vice-versa. Block 9-14 trigger changes in MV2, being only B9 the block that turns ON MV2, while the others turn it OFF. Finally, B16 updates the actuator signals from values in local variables.

---

```

1 (*Read from sensors*)
2 'B1:' MV2.ZSC := DigitalInput(1);

```

---

```

3 MV2.ZSO := DigitalInput(2);
4 L2.Level := AnalogInput(1);
5 (*Constants*)
6 'B2:' L2min := 800;
7       L2max := 1000;
8       T2 := 7;
9 'B3:' MV2.SR.Enable := 1;
10       MV2.SR.S := L2.Level < L2min;
11       MV2.SR.R := L2.Level > L2max;
12       SETD(MV2.SR);
13       MV2.Auto := MV2.SR.Out;
14 (*Timers*)
15 'B4:' MV2.Close_TM.PRE := T2*1000;
16       MV2.Close_TM.Enable := MV2.Cmd_Close;
17       TONR(MV2.Close_TM);
18       MV2.Open_TM.PRE := T2*1000;
19       MV2.Open_TM.Enable := MV2.Cmd_Open;
20       TONR(MV2.Open_TM);
21 'B5:' IF MV2.ZSC THEN
22       MV2.Status :=1; (*MV2 fully closed*)
23 'B6:' ELSE IF MV2.ZSO THEN
24       MV2.Status :=2; (*MV2 fully open*)
25 'B7:' ELSE
26       MV2.Status :=0; (*MV2 in transition*)
27 'B8:' END IF;
28 'B9:' IF MV2.Auto AND (NOT MV2.FC
29       AND NOT MV2.FO) THEN
30       MV2.Cmd_Close := 0;
31       MV2.Cmd_Open := 1; (*Turn ON MV2*)
32 'B10:' IF MV2.Open_TM.DN AND
33        (NOT MV2.ZSO) THEN
34       MV2.FC := 0;
35       MV2.FO := 1;
36 'B11:' END IF;
37 'B12:' ELSE
38       MV2.Cmd_Close := 1; (*Turn OFF MV2*)
39       MV2.Cmd_Open := 0;
40 'B13:' IF MV2.Close_TM.DN
41        AND (NOT MV2.ZSC) THEN
42       MV2.FC := 1;
43       MV2.FO := 0;
44 'B14:' END IF;
45 'B15:' END IF;
46 (*Write to actuators*)
47 'B16:' DigitalOutput(1) := MV2.Cmd_Close
48        DigitalOutput(2) := MV2.Cmd_Open

```

---

Listing 2: Slice of PLC program controlling MV2.

We developed Algorithm 1 to analyze PLC source code, and automatically produce the set of control conditions that trigger changes in the actuator under analysis (MV2; e.g., Listing 2). To design Algorithm 1 we leveraged static program analysis concepts, like control flow analysis, symbolic execution and taint analysis. As a descriptive example, let us apply the algorithm to deduce the conditions that make MV2 to open.

To open MV2 the DigitalOutput(2) must be set to 1 (line 45). First, the algorithm builds the Control Flow Graph (CFG) depicted in Fig. 10. Nodes represent basic blocks<sup>2</sup> and edges show the execution order. Then, the algorithm finds the definitions, two in this case for MV2.Cmd\_Open (lines 30 and 37), but only the path B9-B15-B16 matches the desired output (see Fig. 10). The symbolic execution engine will produce two conditions that satisfy the path. (1) MV2.Auto=1 ∧ MV2.TON\_Open\_TM.DN=0, and

2. A basic block in static program analysis refers to a group of statements that run in sequence (without branches).

(2)  $MV2.Auto=1 \wedge MV2.TON\_Open\_TM.DN=1 \wedge MV2.ZSO=1$ . The dependency analysis shows a link between lines 20 and 31.  $MV2.Open\_TM.DN=1$  is equivalent to ' $MV2.Cmd\_Open=1$  for at least  $T2$  seconds' in B4 or ( $\tau_2 \geq T_2$ ,  $MV2 = T\uparrow$ ). Similarly,  $MV2.Auto$  is defined in line 13, then after processing the latch function we get that  $MV2.Auto$  is equivalent to  $L2.Level < L2min$ .

Algorithm 1 refers to these conditions as symbolic path conditions and they are stored in the R hashmap (line 10, Algorithm 1).

First two expressions in equation (17) (next subsection) describe the symbolic path conditions detailed above. The rest of the control conditions can be automatically extracted from the PLC source code using the Algorithm 1.

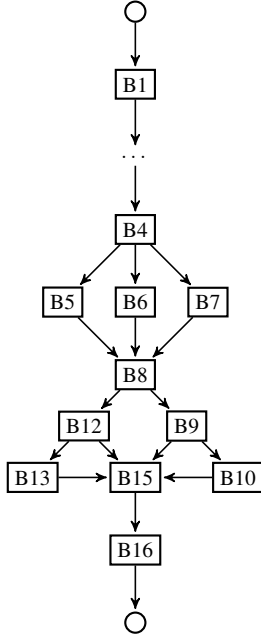


Figure 10: CFG of PLC code controlling MV2 (Listing 2).

Our analysis can be applied to various industrial control processes, as our symbolic path conditions are generic and applicable to all PLC programs written in the “structured text” programming language.

## B. Numerical Values for the Model

In this section, we provide the numerical values of the different parameters used in the modeling and the analysis of the CPS use-case. That is, using the extraction algorithms, we conclude that  $(L1min, L1max) := (500, 800)$ ,  $(L2min, L2max) := (800, 1000)$ ,  $(L3min, L3max) := (800, 1000)$ , and  $(T_1, T_2, T_3) := (9, 7, 7)$ . Furthermore, the numerical values of the rate of change of the water levels L1, L2, and L3; namely,  $F_{L1}$ ,  $F_{L2}$ , and  $F_{L3}$ , have been computed using linear regression.

### 1) In the absence of attacks and when P3 is always ON.

- $F_{L3} \equiv F_{L3}(MV3)$  with  $F_{L3}(0) = -0.15$ ,  $F_{L3}(1) = 0.16$ ,  $F_{L3}(2) = F_{L3}(3) = 0.11$ .
- $F_{L2} \equiv F_{L2}(MV2, P2)$  with  $F_{L2}(0, 0) = 0$ ,  $F_{L2}(1, 0) = 0.46$ ,  $F_{L2}(1, 1) = 0.13$ ,

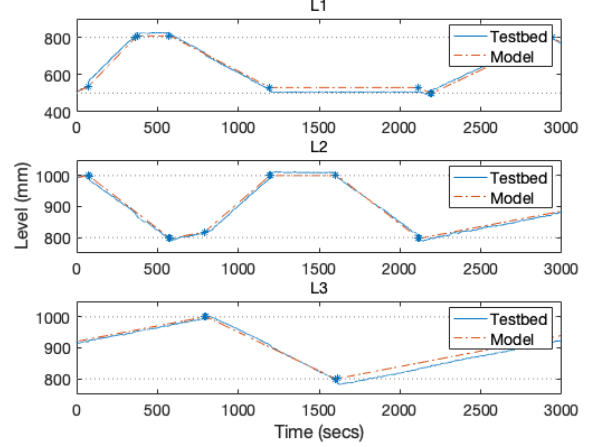


Figure 11: Nominal behavior of the system. Stars show jumps due to actuator transitions.

$$F_{L2}(2, 1) = F_{L2}(3, 1) = -0.28,$$

$$F_{L2}(2, 0) = F_{L2}(3, 0) = 0.29.$$

- $F_{L1} \equiv F_{L1}(MV1, P1)$  with

$$F_{L1}(0, 1) = -0.45, F_{L1}(1, 1) = 0.41,$$

$$F_{L1}(2, 1) = F_{L1}(3, 1) = -0.15,$$

$$F_{L1}(0, 0) = 0, F_{L1}(1, 0) = 0.9,$$

$$F_{L1}(2, 0) = F_{L1}(3, 0) = 0.81.$$

### 2) When only $P2^a$ is attacked and P3 is always ON.

- $F_{L3} \equiv F_{L3}(MV3, P2^a)$  with

$$F_{L3}(MV3, 0) = -0.15 \quad \forall MV3 \in \{0, 1, 2, 3\},$$

and  $F_{L3}(MV3, 1) \equiv F_{L3}(MV3)$  as in 1).

- $F_{L2} \equiv F_{L2}(MV2, MV3, P2^a)$  with

$$F_{L2}(MV2, 1, P2^a) \equiv F_{L2}(MV2, P2^a),$$

$$F_{L2}(MV2, MV3, 0) \equiv F_{L2}(MV2, 0)$$

$$F_{L2}(MV2, 0, P2^a) \equiv F_{L2}(MV2, 0)$$

$$F_{L2}(MV2, 2, 1) \equiv F_{L2}(MV2, 1)$$

as in 1). Furthermore,

$$F_{L2}(MV2, 2, P2^a) \equiv F_{L2}(MV2, 3, P2^a).$$

- $F_{L1} \equiv F_{L1}(MV1, P1)$  as in 1).

### 3) In the absence of attacks and P3 is not always ON.

- $F_{L3} \equiv F_{L3}(MV3, P3)$  with  $F_{L3}(MV3, 1)$  as in 1) and  $F_{L3}(0, 0) = 0$ ,  $F_{L3}(1, 0) = 0.36$ ,  $F_{L3}(2, 0) = F_{L3}(3, 0) = 0.3$ .
- $F_{L2} \equiv F_{L2}(MV2, P2)$  as in 1).
- $F_{L1} \equiv F_{L1}(MV1, P1)$  as in 1).

### 4) When only $P2^a$ is attacked and P3 is not always ON.

- $F_{L3} \equiv F_{L3}(MV3, P3, P2^a)$  with

$$F_{L3}(MV3, P3, 1) \equiv F_{L3}(MV3, P3)$$

$$F_{L3}(MV3, P3, 0) \equiv F_{L3}(0, P3)$$

as in 3).

- $F_{L2} \equiv F_{L2}(MV2, MV3, P2^a)$  as in 2).
- $F_{L1} \equiv F_{L1}(MV1, P1)$  as in 1).

5) **When  $P1^a$  in Attacked and  $P3$  is not Always ON.**

- $F_{L3} \equiv F_{L3}(MV3, P3)$  as in 3).
- $F_{L2} \equiv F_{L2}(MV2, P2, P1^a)$  with

$$F_{L2}(MV2, P2, 1) \equiv F_{L2}(MV2, P2)$$

$$F_{L2}(MV2, P2, 0) \equiv F_{L2}(0, P2)$$

as in 1).

- $F_{L1} \equiv F_{L1}(MV1, MV2, P1^a)$  with

$$F_{L1}(MV1, MV2, 0) =$$

$$F_{L1}(MV1, 0, P1^a) \equiv F_{L1}(MV1, 0)$$

$$F_{L1}(MV1, 1, 1) = F_{L1}(MV1, 2, 1) \equiv F_{L1}(MV1, 1)$$

as in 1).

To test if our model follows the dynamics of the real-world system, we implemented our equations in Matlab's HyEq toolbox [52]. Figure 11 and Figure 12 shows traces of the real-world system labeled as **Testbed** of 3000 seconds under nominal conditions. It also shows the execution of the **Model** during the same period. The model follows the behavior of the Testbed, including triggering actuator transitions synchronously. The different initial conditions of the state vector  $x$  are:

$$\begin{aligned} L1 &= 506.6; & \tau1 &= 1; & MV1 &= \text{ON}; & P1 &= \text{ON}; \\ L2 &= 992.7; & \tau2 &= 1; & MV2 &= \text{ON}; & P2 &= \text{ON}; \\ L3 &= 920.8; & \tau3 &= 1; & MV3 &= \text{ON}; & P3 &= \text{ON}. \end{aligned}$$

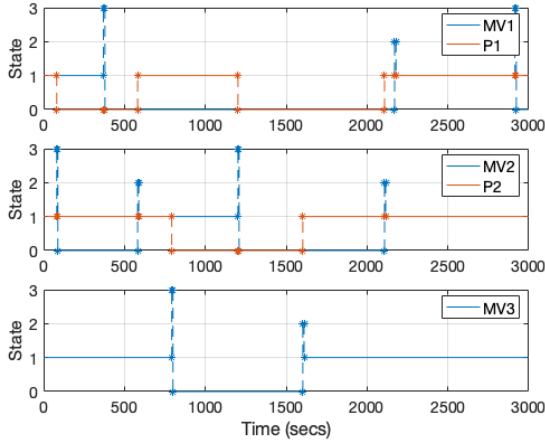


Figure 12: Nominal behavior of the system. Stars show jumps due to actuator transitions.

Figure 13 shows the behavior of our model under multiple operational scenarios of the SWaT testbed, as another example of the fidelity of our model.

### C. New Uniform-Safety Theorems

Our previous results (Theorem 5 below) provide sufficient conditions in terms of infinitesimal inequalities — namely, without using information about solutions to the hybrid system — to guarantee that the set  $K_e$ , on which the barrier function is nonpositive, is forward invariant; namely, the solutions starting from  $K_e$  remain in  $K_e$ . More precisely, under mild conditions on the data  $(C, F, D, G)$  of the hybrid system, we present conditions for which a barrier function guarantees forward invariance

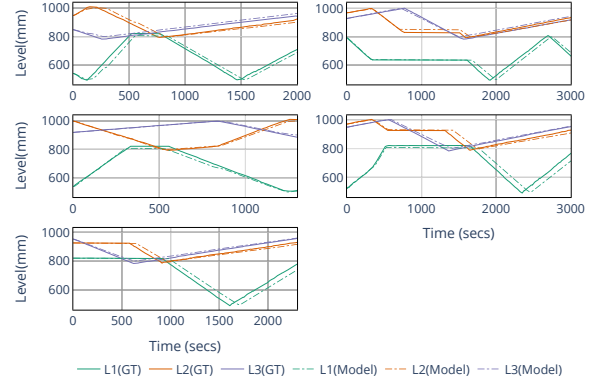


Figure 13: Comparison of multiple traces of the model against ground truth (GT)

of  $K_e$ . It should be noted that the proposed conditions require the barrier function to have, at points where flows are possible, a nonpositive derivative on a neighborhood of the set  $K_e$  and, after a jump from points where jumps are allowed, a nonpositive value.

**Definition 6 (Contingent Cone).** For a set  $K \subset \mathbb{R}^{m_x}$ , according to [7], the *contingent cone* of  $K$  at  $x$  is given by

$$T_K(x) := \left\{ v \in \mathbb{R}^n : \liminf_{h \rightarrow 0^+} \frac{|x + hv|_K}{h} = 0 \right\}. \quad (44)$$

**Theorem 5.** Given a hybrid system  $\mathcal{H} = (C, F, D, G)$  as in (3), suppose that  $F$  is continuous and that there exists a  $C^1$  barrier function candidate  $B$  with respect to  $(X_o, X_u)$  as in (4). The hybrid system  $\mathcal{H}$  is safe with respect to  $(X_o, X_u)$  if (7) and (8) hold and

$$\begin{aligned} \langle \nabla B(x), F(x) \rangle &\leq 0 \quad \forall x \in (U(\partial K_e) \setminus K_e) \cap \\ C: F(x) &\in T_C(x). \end{aligned} \quad (45)$$

□

Note that (45) is a relaxation of (6) in which we constrain the dynamics  $F$  only when it point towards the set  $C$ .

*Proof.* We prove Theorem 1 using Theorem 5. To do so, given an exogenous signal  $u : \mathbb{R}_{\geq 0} \rightarrow \mathcal{U}$ , we introduce the following augmented system

$$\mathcal{H}_u : \begin{cases} \begin{bmatrix} t \\ \dot{x} \end{bmatrix} = \begin{bmatrix} 1 \\ F(x, u(t)) \end{bmatrix} & (t, x) \in C_a \\ \begin{bmatrix} t^+ \\ x^+ \end{bmatrix} = \begin{bmatrix} t \\ G(x, u(t)) \end{bmatrix} & (t, x) \in D_a, \end{cases} \quad (46)$$

where  $C_a := \{(t, x) \in \mathbb{R}_{\geq 0} \times X : (x, u(t)) \in C\}$  and  $D_a := \{(t, x) \in \mathbb{R}_{\geq 0} \times X : (x, u(t)) \in D\}$ . Furthermore, we consider the augmented initial and unsafe sets  $X_{oa} := \mathbb{R}_{\geq 0} \times X_o$  and  $X_{ou} := \mathbb{R}_{\geq 0} \times X_u$ , respectively. Finally, we consider the barrier function candidate  $B_a : \mathbb{R}_{\geq 0} \times X \rightarrow \mathbb{R}$

given by  $B_a(t, x) := B(x)$ . According to Theorem 5, the system (46) is safe with respect  $(X_{oa}, X_{ua})$  if

$$\begin{aligned} \langle \nabla B_a(t, x), F(x, u(t)) \rangle &\leq 0 \quad \forall (t, x) \in \\ &\quad (U(\partial K_{ea}) \setminus K_{ea}) \cap C_a \\ &\quad \text{if } (1, F(x, u(t))) \in T_{C_a}(t, x), \end{aligned} \quad (47)$$

$$B_a(t, G(x, u(t))) \leq 0 \quad \forall (t, x) \in D_a \cap K_{ea}, \quad (48)$$

$$G(x, u(t)) \subset C_a \cup D_a \quad \forall (t, x) \in D_a \cap K_{ea}, \quad (49)$$

where  $K_{ea} := \mathbb{R}_{\geq 0} \times K_e$ . Note that when  $(t, x) \in D_a \cap K_{ea}$ , it follows that  $(x, u(t)) \in D \cap K_u$ ; hence, (48) and (49) follow under (11) and (12), respectively. Furthermore, when  $(t, x) \in (U(\partial K_{ea}) \setminus K_{ea}) \cap C_a$ , it follows that  $(x, u(t)) \in (U(\partial K_u) \setminus K_u) \cap C$ . Furthermore, having  $(1, F(x, u(t))) \in T_{C_a}(t, x)$  implies that there exist  $\{h_i\}_{i \in \mathbb{N}} \subset \mathbb{R}_{\geq 0}$  and  $\{v_i\}_{i \in \mathbb{N}} \subset \mathbb{R}^{m_x+1}$ , with  $v_i := (v_{1i}, v_{2i}) \in \mathbb{R} \times \mathbb{R}^{m_x}$ , such that  $\lim_{i \rightarrow \infty} h_i \rightarrow 0^+$ ,  $\lim_{i \rightarrow \infty} v_i = (1, F(x, u(t)))$ , and  $(t, x) + h_i v_i \in C_a$ . Note that, for each  $i \in \mathbb{N}$ ,

$$x + h_i v_{2i} \in \{x : \exists t \geq 0 : (t, x) \in C_a\}.$$

Hence,  $x + h_i v_{2i} \in C_u$  for all  $i \in \mathbb{N}$ , therefore  $F(x, u(t)) \in T_{C_u}(x)$  and (47) is guaranteed under (10). ■

**Definition 7 (Barrier function certificate for uniform safety).** A  $\mathcal{C}^1$  barrier function candidate with respect to  $(X_o, X_u)$  is a **barrier certificate for uniform safety** with respect to  $(X_o, X_u)$  if (10)-(12) are satisfied. •

## D. Proof of Safety for SWaT Without Attacks

For each  $i \in \{1, 2, 3\}$ , we construct an explicit barrier certificate  $B$  guaranteeing safety of Stage  $i$  with respect to  $(X_{oi}, X_{ui})$  uniformly in  $u_i \in \mathcal{U}_i$ , with  $X_{ui} \subset X_i \setminus X_{si}$ .

We now use Theorem 5 to show safety of  $\mathcal{H}_3$  in (25).

**Theorem 6.** Consider the hybrid system  $\mathcal{H}_3$  in (25). Consider the initial set  $X_{o3}$  in (29) and an unsafe set  $X_{u3} \subset X_3 \setminus X_{s3}$  with  $X_{s3}$  introduced in (35). Assume that there exists  $\sigma > 0$  such that

$$4T_3(F_{L3}(3) + \sigma) < \delta. \quad (50)$$

Then, the hybrid system  $\mathcal{H}_3$  is safe with respect to  $(X_{o3}, X_{u3})$ , and admits the barrier function certificate

$$B(x_3) := (L3 - L3\min)(L3 - L3\max - f(\tau_3, MV3)), \quad (51)$$

where

$$f(\tau_3, MV3) := (F_{L3}(3) + \sigma)[\tau_3 + T_3 * w_f(MV3)],$$

$$w_f(2) := 0, \quad w_f(1) := 1, \quad w_f(3) := 2, \quad w_f(0) := 3.$$

□

*Proof.* Consider the hybrid system  $\mathcal{H}_3$  in (25), the sets  $(X_{o3}, X_{u3})$ , and the scalar function in (51). Note that

$$f(\tau_3, MV3) \in [0, (F_{L3}(3) + \sigma) * 4 * T_3]$$

for all  $(\tau_3, MV3) \in [0, T_3] \times \{0, 1, 2, 3\}$ . Hence, for  $\sigma > 0$  satisfying (50), we conclude that

$$f(\tau_3, MV3) \in [0, \delta]$$

for all  $(\tau_3, MV3) \in [0, T_3] \times \{0, 1, 2, 3\}$ . Thus, (4) is satisfied. Furthermore, let the set

$$K_e := \{x_3 \in X_3 : L3 \in [L3\min, L3\max + f(\tau_3, MV3)]\}.$$

To conclude safety of the third stage using Theorem 5, we start verifying (7) and (8). To do so, we start noting that the set  $K_e \cap D_3$  satisfies

$$K_e \cap D_3 = A_1 \cup A_2,$$

where

$$A_1 := \{x_3 \in X_3 : L3 \in [L3\max, L3\max + f(\tau_3, 1)],$$

$$MV3 = \{1\}\},$$

$$A_2 := \{x_3 \in X_3 : L3 \in [L3\min, L3\max + f(T_3, MV3)],$$

$$\tau_3 = T_3, MV3 \in \{2, 3\}\}.$$

Furthermore, for each  $x_3 \in A_1$ , we have

$$G_3(x_3) = [L3 \ 0 \ 3]^\top$$

Hence,

$$B(G_3(x_3)) = (L3 - L3\min)(L3 - L3\max - f(0, 3)) \leq 0.$$

The latter inequality is true since

$$f(\tau_3, 1) \leq f(0, 3) \quad \forall \tau_3 \in [0, T_3].$$

Similarly, for each  $x_3 \in A_2$ , we have

$$G_3(x_3) = [L3 \ 0 \ \alpha(MV3)]^\top,$$

where  $\alpha(3) := 0$  and  $\alpha(2) := 1$ . Hence,

$$B(G_3(x_3)) =$$

$$(L3 - L3\min)(L3 - L3\max - f(0, \alpha(MV3))) \leq 0.$$

The latter inequality is true since

$$f(\tau_3, 3) \leq f(0, 0) \quad \forall \tau_3 \in [0, T_3], \text{ and}$$

$$f(\tau_3, 2) \leq f(0, 1) \quad \forall \tau_3 \in [0, T_3].$$

Hence, we conclude that (7) is satisfied. Moreover, to show (8), we note that  $C_3 \cup D_3 = X_3$  and  $G_3(x_3) \in X_3$  for all  $x_3 \in D_3$ .

Next, to verify (6), we start noting that the set  $U(\partial K_e) \setminus K_e$  is give by

$$U(\partial K_e) \setminus K_e = B_1 \cup B_2,$$

where, for some  $\epsilon > 0$ ,

$$B_1 := \{x_3 \in X_3 : L3 \in (L3\min - \epsilon, L3\min)\},$$

and

$$B_2 := \{x_3 \in X_3 : L3 \in (L3\max + f(\tau_3, MV3))[1, \epsilon)\}.$$

Furthermore, the set  $C_3$  can be explicitly expressed as

$$C_3 := \text{cl}(X_3 \setminus D_3) = \bigcap_{i=1}^3 \bar{D}_{3i},$$

where

$$\bar{D}_{31} := \{x_3 \in X_3 : L3 \geq L3\min \cup MV3 \in \{1, 2, 3\}\},$$

$$\bar{D}_{32} := \{x_3 \in X_3 : \tau_3 \leq T_3 \cup MV3 \in \{0, 1\}\},$$

$$\bar{D}_{33} := \{x_3 \in X_3 : L3 \leq L3\max \cup MV3 \in \{0, 2, 3\}\}.$$

This is equivalent to

$$C_3 := \bigcup_{i=1}^4 C_{3i},$$

where

$$\begin{aligned} C_{31} &:= \{x_3 \in X_3 : L3_{\min} \leq L3 \leq L3_{\max}\}, \\ C_{32} &:= \{x_3 \in X_3 : L3 \leq L3_{\max}, MV3 \in \{1, 2, 3\}\}, \\ C_{33} &:= \{x_3 \in X_3 : L3 \geq L3_{\min}, MV3 \in \{0, 2, 3\}\}, \\ C_{34} &:= \{x_3 \in X_3 : MV3 \in \{2, 3\}\}. \end{aligned}$$

Hence, we conclude that

$$\begin{aligned} (U(\partial K_e) \setminus K_e) \cap C_3 = \\ (B_1 \cap C_{32}) \cup (B_2 \cap C_{33}) \cup (B_1 \cap C_{34}) \cup (B_2 \cap C_{34}), \end{aligned}$$

with

$$\begin{aligned} B_1 \cap C_{32} &= \{x_3 \in X_3 : \\ L3 &\in (L3_{\min} - \epsilon, L3_{\min}), MV3 \in \{1, 2, 3\}\}, \end{aligned}$$

$$\begin{aligned} B_2 \cap C_{33} &= \{x_3 \in X_3 : \\ L3 &\in (L3_{\max} + f(\tau_3, MV3), L3_{\max} + f(\tau_3, MV3) + \epsilon), \\ MV3 &\in \{0, 2, 3\}\}, \end{aligned}$$

$$\begin{aligned} B_1 \cap C_{34} &= \{x_3 \in X_3 : \\ L3 &\in (L3_{\min} - \epsilon, L3_{\min}), MV3 \in \{2, 3\}\}, \end{aligned}$$

$$\begin{aligned} B_2 \cap C_{34} &= \{x_3 \in X_3 : \\ L3 &\in (L3_{\max} + f(\tau_3, MV3), L3_{\max} + f(\tau_3, MV3) + \epsilon), \\ MV3 &\in \{2, 3\}\}. \end{aligned}$$

Next, we evaluate the scalar product  $\langle \nabla B(x_3), F_3(x_3) \rangle$  at each  $x_3 \in (U(\partial K_e) \setminus K_e) \cap C_3$ . Indeed, note that

$$\nabla B(x_3) = \begin{bmatrix} 2L3 - (L3_{\max} + L3_{\min}) - f(\tau_3, MV3) \\ -(F_{L3}(3) + \sigma)(L3 - L3_{\min}) \\ \star \end{bmatrix}.$$

Hence,

$$\begin{aligned} \langle \nabla B(x_3), F(x_3) \rangle &= F_{L3}(MV3)(L3 - L3_{\max} - \\ &\quad f(\tau_3, MV3)) + (L3 - L3_{\min}) \\ &\quad (F_{L3}(MV3) - F_{L3}(3) - \sigma). \end{aligned}$$

Next, we distinguish the following two situations:

- 1) When  $x_3 \in (B_1 \cap C_{32}) \cup (B_1 \cap C_{34})$ . In this case, we conclude that  $F_{L3}(MV3) \geq F_{L3}(3)$ ,  $|L3 - L3_{\min}| \leq \epsilon$ , and  $L3_{\max} - L3 > L3_{\max} - L3_{\min}$ . Hence,
$$\begin{aligned} \langle \nabla B(x_3), F_3(x_3) \rangle &\leq F_{L3}(MV3)(L3 - L3_{\max}) + \\ &\quad \epsilon |F_{L3}(MV3) - (F_{L3}(3) + \sigma)| \\ &\leq -F_{L3}(3) * (L3_{\max} - L3_{\min}) + \\ &\quad \epsilon |F_{L3}(MV3) - (F_{L3}(3) + \sigma)|. \end{aligned}$$

Hence, for  $\epsilon$  sufficiently small, we conclude that

$$\langle \nabla B(x_3), F_3(x_3) \rangle \leq 0.$$

- 2) When  $x_3 \in (B_2 \cap C_{33}) \cup (B_2 \cap C_{34})$ . In this case, we conclude that  $|L3 - L3_{\max} - f(\tau_3, MV3)| \leq \epsilon$ ,  $L3 - L3_{\min} \geq L3_{\max} - L3_{\min}$ , and

$$F_{L3}(MV3) - (F_{L3}(3) + \sigma) \leq -\sigma.$$

Hence,

$$\begin{aligned} \langle \nabla B(x_3), F_3(x_3) \rangle &\leq |F_{L3}(MV3)|\epsilon + \\ &\quad (L3 - L3_{\min})[F_{L3}(MV3) - (F_{L3}(3) + \sigma)] \\ &\leq |F_{L3}(MV3)| * \epsilon - \sigma * (L3_{\max} - L3_{\min}). \end{aligned}$$

Hence, for  $\epsilon$  sufficiently small, we conclude that

$$\langle \nabla B(x_3), F_3(x_3) \rangle \leq 0. \quad \blacksquare$$

Next, using Theorem 1, we show safety of  $\mathcal{H}_2$  in (26) uniformly in  $u_2 \in \mathcal{U}_2$ .

**Theorem 7.** Consider the hybrid system  $\mathcal{H}_2$  in (26). Consider the initial set  $X_{o2}$  in (30) and the unsafe set  $X_{u2} \subset X_2 \setminus X_{s2}$  with  $X_{s2}$  introduced in (34). Assume that there exist  $\sigma_h > 0$  and  $\sigma_g > 0$  such that

$$4T_2(F_{L2}(2, 0) + \sigma_h) < \delta, \quad (52)$$

$$4T_2(F_{L2}(2, 1) + \sigma_g) < \delta. \quad (53)$$

Then, the system  $\mathcal{H}_2$  in (26) is safe with respect to  $(X_{o2}, X_{u2})$  uniformly in  $u_2 \in \mathcal{U}_2$ . Moreover, the system  $\mathcal{H}_2$  admits a barrier function certificate for uniform safety given by

$$\begin{aligned} B(x_2) &:= (L2 - L2_{\min} + g(\tau_2, MV2)) \times \\ &\quad (L2 - L2_{\max} - h(\tau_2, MV2)), \end{aligned} \quad (54)$$

where

$$g(\tau_2, MV2) := (-F_{L2}(2, 1) + \sigma_g) * [\tau_2 + T_2 * w_g(MV2)],$$

$$w_g(3) := 0, \quad w_g(0) := 1, \quad w_g(2) := 2, \quad w_g(1) := 3,$$

and

$$h(\tau_2, MV2) := (F_{L2}(2, 0) + \sigma_h) * [\tau_2 + T_2 * w_h(MV2)],$$

$$w_h(2) := 0, \quad w_h(1) := 1, \quad w_h(3) := 2, \quad w_h(0) := 3. \quad \square$$

*Proof.* Consider the hybrid system  $\mathcal{H}_2$  in (26) and the scalar function in (51). Note that, for all  $(\tau_2, MV2) \in [0, T_2] \times \{0, 1, 2, 3\}$ , we have

$$g(\tau_2, MV2) \in [0, ([-F_{L2}(2, 1) + \sigma_g] * T_2 * 4)],$$

and

$$h(\tau_2, MV2) \in [0, ([F_{L2}(2, 0) + \sigma_h] * T_2 * 4)].$$

Hence, for  $\sigma_g$  and  $\sigma_h$  satisfying (82) and (83), we conclude that, for all  $(\tau_2, MV2) \in [0, T_2] \times \{0, 1, 2, 3\}$

$$g(\tau_2, MV2) \in [0, \delta] \text{ and } h(\tau_2, MV2) \in [0, \delta].$$

Thus, (4) is satisfied. Furthermore, note that

$$\begin{aligned} K_w &= \{x_2 \in X_2 : \\ L2 &\in [L2_{\min} - g(\tau_2, MV2), L2_{\max} + h(\tau_2, MV2)]\} \times \mathcal{U}_2. \end{aligned}$$

To apply Theorem 1, we start verifying (11) and (12). To do so, we start noting that the set  $K_w \cap D_2$  satisfies

$$K_w \cap D_2 = A_1 \cup A_2 \cup A_3 \cup A_4,$$

To verify (10), we have the following two situations:

- 1) When  $(x_2, u_2) \in (C_{21} \cap B_1) \cup (C_{22} \cap B_1) \cup (C_{25} \cap B_1) \cup (C_{26} \cap B_1)$ , we conclude that

$$\begin{aligned} \langle \nabla B(x_2), F_2(x_2) \rangle &\leq \sigma_g [L_2 - L_{2\max} - h(\tau_2, \text{MV2})] + \\ &\epsilon [F_{L_2}(2, 0) + \sigma_h] \leq \sigma_g [L_{2\min} - L_{2\max} - \\ &h(\tau_2, \text{MV2}) - g(\tau_2, \text{MV2})] + \epsilon [F_{L_2}(2, 0) + \sigma_h] \end{aligned}$$

Finally, for  $\epsilon$  sufficiently small, we conclude that

$$\langle \nabla B(x_2), F_2(x_2) \rangle \leq 0.$$

- 2) When  $(x_2, u_2) \in (C_{25} \cap B_2) \cup (C_{26} \cap B_2) \cup (C_{27} \cap B_2) \cup (C_{28} \cap B_2)$ , we conclude that

$$\begin{aligned} \langle \nabla B(x_2), F_2(x_2) \rangle &\leq \\ &-\sigma_h [L_2 - L_{2\min} + g(\tau_2, \tau_3)] + \\ &\epsilon [-F_{L_2}(2, 1) + \sigma_g + F_{L_2}(\text{MV2}, \text{P2})] \leq \\ &-\sigma_h [L_{2\max} - L_{2\min} + g(\tau_2, \tau_3) + h(\tau_2, \tau_3)] + \\ &\epsilon [-F_{L_2}(2, 1) + \sigma_g + F_{L_2}(\text{MV2}, \text{P2})]. \end{aligned}$$

Hence, for  $\epsilon$  sufficiently small, we conclude that

$$\langle \nabla B(x_2), F_2(x_2) \rangle \leq 0. \quad \blacksquare$$

The same statement as the one in Theorem 7 can be formulated for Stage 1, mutatis mutandis.

**Theorem 8.** Consider the hybrid system  $\mathcal{H}_1$  in (27). Consider the initial set  $X_{o1}$  in (29) and the unsafe set  $X_{u1} \subset X_1 \setminus X_{s1}$  with  $X_{s1}$  introduced in (33). Assume that there exist  $\sigma_h > 0$  and  $\sigma_g > 0$  such that

$$4T_1(F_{L_2}(2, 0) + \sigma_h) < \delta, \quad (55)$$

$$4T_2(F_{L_2}(2, 1) + \sigma_g) < \delta. \quad (56)$$

Then, the system  $\mathcal{H}_1$  in (27) is safe with respect to  $(X_{o1}, X_{u1})$  uniformly in  $u_1 \in \mathcal{U}_1$ . Moreover, the system  $\mathcal{H}_1$  admits a barrier function certificate for uniform safety given by

$$\begin{aligned} B(x_1) := &(L_1 - L_{1\min} + g(\tau_1, \text{MV1})) \times \\ &(L_1 - L_{1\max} - h(\tau_1, \text{MV1})), \end{aligned} \quad (57)$$

where

$$g(\tau_1, \text{MV1}) := (-F_{L_1}(2, 1) + \sigma_g) * [\tau_1 + T_1 * w_g(\text{MV1})],$$

$$w_g(3) := 0, \quad w_g(0) := 1, \quad w_g(2) := 2, \quad w_g(1) := 3,$$

and

$$h(\tau_1, \text{MV1}) := (F_{L_1}(2, 0) + \sigma_h) * [\tau_1 + T_1 * w_h(\text{MV1})],$$

$$w_h(2) := 0, \quad w_h(1) := 1, \quad w_h(3) := 2, \quad w_h(0) := 3. \quad \square$$

The system  $\mathcal{H}_1$  is safe with respect to  $(X_{o1}, X_{u1})$  uniformly in  $u_1 \in \mathcal{U}_1$  when, for some  $\sigma_h > 0$  and  $\sigma_g$ ,

$$4T_1(F_{L_1}(2, 0) + \sigma_h) < \delta, \quad (58)$$

$$4T_1(F_{L_1}(2, 1) + \sigma_g) < \delta. \quad (59)$$

Moreover,  $\mathcal{H}_1$  admits the barrier function certificate

$$\begin{aligned} B(x_1) := &(L_1 - L_{1\min} + g(\tau_1, \text{MV1})) \times \\ &(L_1 - L_{1\max} - h(\tau_1, \text{MV1})). \end{aligned} \quad (60)$$

Finally, the combination of the previous statements allows us to conclude safety for the entire plant.

**Theorem 9.** Consider the hybrid system  $\mathcal{H}$  composed by the cascaded interconnection of  $\mathcal{H}_1$ ,  $\mathcal{H}_2$ , and  $\mathcal{H}_3$ . Consider the initial set  $X_o$  in (28) and an unsafe set

$$X_u \subset (X_1 \setminus X_{s1}) \times (X_2 \setminus X_{s2}) \times (X_3 \setminus X_{s3}),$$

where  $(X_{s1}, X_{s2}, X_{s3})$  are introduced in (33), (34), and (35), respectively. Assume that there exists  $\sigma > 0$ ,  $\sigma_h > 0$ , and  $\sigma_g > 0$  such that (50), (82), (83), (58), and (59) hold. Then, the hybrid system  $\mathcal{H}$  is safe with respect to  $(X_o, X_u)$ .  $\square$

*Proof.* Consider the hybrid system  $\mathcal{H}_2$  in (26) and the scalar function in (51). Note that, for all  $(\tau_2, \text{MV2}) \in [0, T_2] \times \{0, 1, 2, 3\}$ , we have

$$g(\tau_2, \text{MV2}) \in [0, ([-F_{L_2}(2, 1) + \sigma_g] * T_2 * 4)],$$

and

$$h(\tau_2, \text{MV2}) \in [0, ([F_{L_2}(2, 0) + \sigma_h] * T_2 * 4)].$$

Hence, for  $\sigma_g$  and  $\sigma_h$  satisfying (82) and (83), we conclude that, for all  $(\tau_2, \text{MV2}) \in [0, T_2] \times \{0, 1, 2, 3\}$ , we have

$$g(\tau_2, \text{MV2}) \in [0, \delta] \text{ and } h(\tau_2, \text{MV2}) \in [0, \delta].$$

Thus, (4) is satisfied. Furthermore, note that

$$K_w = \{x_2 \in X_2 :$$

$$L_2 \in [L_{2\min} - g(\tau_2, \text{MV2}), L_{2\max} + h(\tau_2, \text{MV2})]\} \times \mathcal{U}_2.$$

To apply Theorem 1, we start verifying (11) and (12). To do so, we start noting that the set  $K_w \cap D_2$  satisfies

$$K_w \cap D_2 = A_1 \cup A_2 \cup A_3 \cup A_4,$$

where

$$A_1 := \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 :$$

$$L_2 \in [L_{2\max}, L_{2\max} + h(\tau_2, 1)], \text{ MV2} = 1\},$$

$$A_2 := \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 :$$

$$L_2 \in [L_{2\min} - g(\tau_2, 0), L_{2\min}], \text{ MV2} = 0\},$$

$$A_3 := \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 :$$

$$L_2 \in [L_{2\min} - g(T_2, \text{MV2}), L_{2\max} + h(T_2, \text{MV2})], \\ \tau_2 = T_2, \text{ MV2} \in \{2, 3\}\},$$

$$A_4 := \{(x_2, u_2) \in (X_2 \times \mathcal{U}_2) \setminus (A_1 \cup A_2 \cup A_3) :$$

$$L_2 \in [L_{2\min} - g(\tau_2, \text{MV2}), L_{2\max} + h(\tau_2, \text{MV2})], \\ \text{P2} = 1, L_3 \geq L_{3\max}\} \cup$$

$$\{(x_2, u_2) \in (X_2 \times \mathcal{U}_2) \setminus (A_1 \cup A_2 \cup A_3) :$$

$$L_2 \in [L_{2\min} - g(\tau_2, \text{MV2}), L_{2\max} + h(\tau_2, \text{MV2})], \\ \text{P2} = 1, \text{ MV3} \in \{0, 3\}\},$$

$$A_5 := \{(x_2, u_2) \in (X_2 \times \mathcal{U}_2) \setminus (A_1 \cup A_2 \cup A_3) :$$

$$L_2 \in [L_{2\min} - g(\tau_2, \text{MV2}), L_{2\max} + h(\tau_2, \text{MV2})], \\ \text{P2} = 0, L_3 \leq L_{3\min}\} \cup$$

$$\{(x_2, u_2) \in (X_2 \times \mathcal{U}_2) \setminus (A_1 \cup A_2 \cup A_3) :$$

$$L_2 \in [L_{2\min} - g(\tau_2, \text{MV2}), L_{2\max} + h(\tau_2, \text{MV2})], \\ \text{P2} = 0, \text{ MV3} \in \{1, 2\}\}.$$



Note that, for each  $(x_2, u_2) \in A_1$ , we have

$$G_2(x_2, u_2) \in [\text{L2} \ 0 \ 3 \ \{0, 1\}]^\top.$$

Hence,

$$B(G_2(x_2, u_2)) = (\text{L2} - \text{L2min} + g(0, 3))(\text{L2} - \text{L2max} - h(0, 3)) \leq 0.$$

The latter inequality is true since  $g(0, 3) \geq 0$ ,

$$\text{L2} \geq \text{L2max} \quad \forall (x_2, u_2) \in A_1,$$

and

$$h(\tau_2, 1) \leq h(0, 3) \quad \forall \tau_2 \in [0, T_2].$$

Similarly, for each  $(x_2, u_2) \in A_2$ , we have

$$G_2(x_2, u_2) = [\text{L2} \ 0 \ 2 \ \{0, 1\}]^\top.$$

Hence,

$$B(G_2(x_2, u_2)) = (\text{L2} - \text{L2min} + g(0, 2))(\text{L2} - \text{L2max} - h(0, 2)) \leq 0.$$

The latter inequality is true since  $h(0, 2) \geq 0$ ,

$$\text{L2} \leq \text{L2min} \quad \forall (x_2, u_2) \in A_2,$$

and

$$g(\tau_2, 0) \leq g(0, 2) \quad \forall \tau_2 \in [0, T_2].$$

Now, for each  $(x_2, u_2) \in A_3$ ,

$$G_2(x_2) \in [\text{L2} \ 0 \ \alpha(\text{MV2}) \ \{0, 1\}]^\top,$$

where  $\alpha(3) := 0$  and  $\alpha(2) := 1$ . Hence,

$$B(G_2(x_2)) = (\text{L2} - \text{L2min} + g(0, \alpha(\text{MV2}))) * (\text{L2} - \text{L2max} - h(0, \alpha(\text{MV2}))) \leq 0.$$

The latter inequality is true since

$$h(7, 3) \leq h(0, 0), \quad h(7, 2) \leq h(0, 1), \\ g(7, 3) \leq g(0, 0), \quad \text{and} \quad g(7, 2) \leq g(0, 1).$$

Finally, for each  $(x_2, u_2) \in A_4 \cup A_5$ , we have

$$G_2(x_2, u_2) \in [\text{L2} \ \tau_2 \ \text{MV2} \ \{0, 1\}]^\top.$$

Hence,

$$B(G_2(x_2, u_2)) = (\text{L2} - \text{L2min} + g(\tau_2, \text{MV2})) * (\text{L2} - \text{L2max} - h(\tau_2, \text{MV2})) \leq 0.$$

The latter inequality is true since  $(A_4 \cup A_5) \subset K_w$ . Hence, we conclude that (11) is satisfied. Moreover, to verify (12), we note that  $C_2 \cup D_2 = X_2 \times \mathcal{U}_2$  and  $G_2(x_2, u_2) \in X_2$  for all  $(x_2, u_2) \in D_2$ .

Next, to verify (10), we start noting that

$$U(\partial K_w) \setminus K_w = B_1 \cup B_2,$$

where, for some  $\epsilon > 0$  sufficiently small,

$$B_1 := \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \text{L2} \in (\text{L2min} - g(\tau_2, \text{MV2}) - \epsilon, \text{L2min} - g(\tau_2, \text{MV2}))\},$$

and

$$B_2 := \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \text{L2} \in (\text{L2max} + h(\tau_2, \text{MV2}), \text{L2max} + h(\tau_2, \text{MV2}) + \epsilon)\}.$$

Furthermore, the set  $C_2$  can be explicitly expressed as

$$C_2 = \bigcup_{i=1}^8 C_{2i},$$

where

$$C_{2i} := C_{2i}^a \cup C_{2i}^b \quad \forall i \in \{2, 4, 6, 8\},$$

and

$$\begin{aligned} C_{21} &:= \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \text{L2} \leq \text{L2max}, \text{MV2} = \{1, 2, 3\}, \\ &\quad \text{P2} = \{0\}, \text{MV3} = \{0, 3\}\}, \\ C_{22}^a &:= \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \text{L2} \leq \text{L2max}, \text{MV2} = \{1, 2, 3\}, \\ &\quad \text{L3} \leq \text{L3max}, \text{P2} = \{0\}, \text{MV3} = \{0, 3\}\}, \\ C_{22}^b &:= \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \text{L2} \leq \text{L2max}, \text{MV2} = \{1, 2, 3\}, \\ &\quad \text{L3} \leq \text{L3max}, \text{P2} = \{1\}, \text{MV3} = \{1, 2\}\}, \\ C_{23} &:= \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \\ &\quad \text{L2min} \leq \text{L2} \leq \text{L2max}, \text{P2} = \{0\}, \text{MV3} = \{0, 3\}\}, \\ C_{24}^a &:= \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \text{L2min} \leq \text{L2} \leq \text{L2max}, \\ &\quad \text{L3} \leq \text{L3max}, \text{P2} = \{0\}, \text{MV3} = \{0, 3\}\}, \\ C_{24}^b &:= \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \text{L2min} \leq \text{L2} \leq \text{L2max}, \\ &\quad \text{L3} \leq \text{L3max}, \text{P2} = \{1\}, \text{MV3} = \{1, 2\}\}, \\ C_{25} &:= \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \\ &\quad \text{MV2} \in \{2, 3\}, \text{P2} = \{0\}, \text{MV3} = \{0, 3\}\}, \\ C_{26}^a &:= \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \text{MV2} \in \{2, 3\}, \text{L3} \leq \text{L3max}, \\ &\quad \text{P2} = \{0\}, \text{MV3} = \{0, 3\}\}, \\ C_{26}^b &:= \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \text{MV2} \in \{2, 3\}, \text{L3} \leq \text{L3max}, \\ &\quad \text{P2} = \{1\}, \text{MV3} = \{1, 2\}\}, \\ C_{27} &:= \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \text{L2} \geq \text{L2min}, \text{MV2} = \{0, 2, 3\}, \\ &\quad \text{P2} = \{0\}, \text{MV3} = \{0, 3\}\}, \\ C_{28}^a &:= \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \text{L2} \geq \text{L2min}, \text{MV2} = \{0, 2, 3\}, \\ &\quad \text{L3} \in [\text{L3min}, \text{L3max}], \text{P2} = \{0\}, \text{MV3} = \{0, 3\}\}, \\ C_{28}^b &:= \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \text{L2} \geq \text{L2min}, \text{MV2} = \{0, 2, 3\}, \\ &\quad \text{L3} \in [\text{L3min}, \text{L3max}], \text{P2} = \{1\}, \text{MV3} = \{1, 2\}\}. \end{aligned}$$

Hence, we conclude that

$$\begin{aligned} (U(\partial K_w) \setminus K_w) \cap C_2 &= (C_{21} \cap B_1) \cup (C_{22} \cap B_1) \cup \\ &= (C_{25} \cap B_1) \cup (C_{25} \cap B_2) \cup (C_{26} \cap B_1) \cup (C_{26} \cap B_2) \cup \\ &= (C_{27} \cap B_2) \cup (C_{28} \cap B_2), \end{aligned}$$

with

$$\begin{aligned} C_{21} \cap B_1 &= \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \\ &\quad \text{L2} \in (\text{L2min} - g(\tau_2, \text{MV2}) - \epsilon, \text{L2min} - g(\tau_2, \text{MV2})), \\ &\quad \text{MV2} = \{1, 2, 3\}, \text{P2} = \{0\}, \text{MV3} = \{0, 3\}\}, \end{aligned}$$

$$\begin{aligned} C_{22} \cap B_1 &= \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \\ &\quad \text{L2} \in (\text{L2min} - g(\tau_2, \text{MV2}) - \epsilon, \text{L2min} - g(\tau_2, \text{MV2})), \\ &\quad \text{MV2} = \{1, 2, 3\}, \text{L3} \leq \text{L3max}, \text{P2} = \{0\}, \\ &\quad \text{MV3} = \{0, 3\}\} \cup \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \\ &\quad \text{L2} \in (\text{L2min} - g(\tau_2, \text{MV2}) - \epsilon, \text{L2min} - g(\tau_2, \text{MV2})), \\ &\quad \text{MV2} = \{1, 2, 3\}, \text{L3} \leq \text{L3max}, \text{P2} = \{1\}, \text{MV3} = \{1, 2\}\}, \end{aligned}$$

$$C_{25} \cap B_1 = \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \\ L2 \in (L2_{\min} - g(\tau_2, MV2) - \epsilon, L2_{\min} - g(\tau_2, MV2)), \\ MV2 \in \{2, 3\}, P2 = \{0\}, MV3 = \{0, 3\}\},$$

$$C_{26} \cap B_1 = \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \\ L2 \in (L2_{\min} - g(\tau_2, MV2) - \epsilon, L2_{\min} - g(\tau_2, MV2)), \\ MV2 \in \{2, 3\}, L3 \leq L3_{\max}, P2 = \{0\}, MV3 = \{0, 3\}\} \cup \\ \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \\ L2 \in (L2_{\min} - g(\tau_2, MV2) - \epsilon, L2_{\min} - g(\tau_2, MV2)), \\ MV2 \in \{2, 3\}, L3 \leq L3_{\max}, P2 = \{1\}, MV3 = \{1, 2\}\},$$

$$C_{25} \cap B_2 = \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \\ L2 \in (L2_{\max} + h(\tau_2, MV2), L2_{\max} + h(\tau_2, MV2) + \epsilon), \\ MV2 \in \{2, 3\}, P2 = \{0\}, MV3 = \{0, 3\}\},$$

$$C_{26} \cap B_2 = \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \\ L2 \in (L2_{\max} + h(\tau_2, MV2), L2_{\max} + h(\tau_2, MV2) + \epsilon), \\ MV2 \in \{2, 3\}, L3 \leq L3_{\max}, P2 = \{0\}, MV3 = \{0, 3\}\} \cup \\ \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \\ L2 \in (L2_{\max} + h(\tau_2, MV2), L2_{\max} + h(\tau_2, MV2) + \epsilon), \\ MV2 \in \{2, 3\}, L3 \leq L3_{\max}, P2 = \{1\}, MV3 = \{1, 2\}\},$$

$$C_{27} \cap B_2 = \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \\ L2 \in (L2_{\max} + h(\tau_2, MV2), L2_{\max} + h(\tau_2, MV2) + \epsilon), \\ MV2 = \{0, 2, 3\}, P2 = \{0\}, MV3 = \{0, 3\}\},$$

$$C_{28} \cap B_2 = \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \\ L2 \in (L2_{\max} + h(\tau_2, MV2), L2_{\max} + h(\tau_2, MV2) + \epsilon), \\ MV2 = \{0, 2, 3\}, L3 \in [L3_{\min}, L3_{\max}], P2 = \{1\}, \\ MV3 = \{1, 2\}\} \cup \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : \\ L2 \in (L2_{\max} + h(\tau_2, MV2), L2_{\max} + h(\tau_2, MV2) + \epsilon), \\ MV2 = \{0, 2, 3\}, L3 \in [L3_{\min}, L3_{\max}], P2 = \{0\}, \\ MV3 = \{0, 3\}\},$$

Next, we evaluate the scalar product  $\langle \nabla B(x_2), F_2(x_2) \rangle$  for each  $(x_2, u_2) \in (U(\partial K_w) \setminus K_w) \cap C_2$ . To do so, we note that

$$\nabla B(x_2) = [\nabla_{L2} B(x_2) \quad \nabla_{\tau_2} B(x_2) \quad \star \quad \star]^\top,$$

where

$$\nabla_{L2} B(x_2) := 2L2 - (L2_{\max} + L2_{\min}) - h(\tau_2, MV2) + g(\tau_2, MV2),$$

$$\nabla_{\tau_2} B(x_2) := -[L2 - L2_{\min} + g(\tau_2, \tau_3)] \star \\ [F_{L2}(2, 0) + \sigma_h] + \\ [L2 - L2_{\max} - h(\tau_2, MV2)] \star \\ [-F_{L2}(2, 1) + \sigma_g].$$

Hence,

$$\langle \nabla B(x_2), F(x_2) \rangle = F_{L2}(MV2, P2) \star \\ [2L2 - (L2_{\min} + L2_{\max}) - h(\tau_2, MV2) + \\ g(\tau_2, MV2)] - \\ [L2 - L2_{\min} + g(\tau_2, \tau_3)][F_{L2}(2, 0) + \sigma_h] + \\ [L2 - L2_{\max} - h(\tau_2, MV2)][-F_{L2}(2, 1) + \sigma_g],$$

and here we distinguish the following two situations:

- 1) When  $(x_2, u_2) \in (C_{21} \cap B_1) \cup (C_{22} \cap B_1) \cup (C_{25} \cap B_1) \cup (C_{26} \cap B_1)$ , we conclude that

$$F_{L2}(MV2, P2) \geq F_{L2}(2, 1),$$

$$-\epsilon \leq L2 - L2_{\min} + g(\tau_2, \tau_3) \leq 0.$$

Hence,

$$\langle \nabla B(x_2), F_2(x_2) \rangle \leq F_{L2}(MV2, P2) \times \\ [L2 - L2_{\max} - h(\tau_2, MV2)] + \\ \epsilon[F_{L2}(2, 0) + \sigma_h] + \\ [L2 - L2_{\max} - h(\tau_2, MV2)][-F_{L2}(2, 1) + \sigma_g]$$

and, thus,

$$\langle \nabla B(x_2), F_2(x_2) \rangle \leq \sigma_g [L2 - L2_{\max} - h(\tau_2, MV2)] + \\ \epsilon[F_{L2}(2, 0) + \sigma_h] \leq \\ \sigma_g [L2_{\min} - L2_{\max} - h(\tau_2, MV2) - g(\tau_2, MV2)] + \\ \epsilon[F_{L2}(2, 0) + \sigma_h]$$

Finally, for  $\epsilon$  sufficiently small, we conclude that

$$\langle \nabla B(x_2), F_2(x_2) \rangle \leq 0.$$

- 2) When  $(x_2, u_2) \in (C_{25} \cap B_2) \cup (C_{26} \cap B_2) \cup (C_{27} \cap B_2) \cup (C_{28} \cap B_2)$ , we conclude that

$$F_{L2}(MV2, P2) \leq F_{L2}(2, 0),$$

$$0 \leq L2 - L2_{\max} - h(\tau_2, \tau_3) \leq \epsilon.$$

Hence,

$$\langle \nabla B(x_2), F_2(x_2) \rangle \leq F_{L2}(MV2, P2) \star \\ [L2 - L2_{\min} + g(\tau_2, MV2) + \epsilon] - \\ [L2 - L2_{\min} + g(\tau_2, \tau_3)][F_{L2}(2, 0) + \sigma_h] + \\ \epsilon[-F_{L2}(2, 1) + \sigma_g]$$

and, thus,

$$\langle \nabla B(x_2), F_2(x_2) \rangle \leq \\ -\sigma_h [L2 - L2_{\min} + g(\tau_2, \tau_3)] + \\ \epsilon[-F_{L2}(2, 1) + \sigma_g + F_{L2}(MV2, P2)] \leq \\ -\sigma_h [L2_{\max} - L2_{\min} + g(\tau_2, \tau_3) + h(\tau_2, \tau_3)] + \\ \epsilon[-F_{L2}(2, 1) + \sigma_g + F_{L2}(MV2, P2)].$$

Hence, for  $\epsilon$  sufficiently small, we conclude that

$$\langle \nabla B(x_2), F_2(x_2) \rangle \leq 0.$$

■

## E. Modeling Actuator Attacks in SWaT

In this section, we assume that the attacker has the ability to falsify the control signals ( $MV1$ ,  $MV2$ ,  $MV3$ ) and ( $P1$ ,  $P2$ ,  $P3$ ) that the controllers ( $C1$ ,  $C2$ ,  $C3$ ) send to the corresponding motor valves and pumps. That is, we let  $w_i := (w_{ip}, w_{im}) \in \{0, 1, 2\} \times \{0, 1\}$ , for all  $i \in \{1, 2, 3\}$ , be the signal sent by the attacker to the  $i$ -th pump and to the  $i$ -th motor valve, respectively. In the presence of attacks, the variables ( $Pi$ ,  $MVi$ ) do not necessarily correspond to the actual states of the  $i$ -th motor valve and the  $i$ -th pump, respectively. For this reason, we introduce the extra variables ( $MV1^a, MV2^a, MV3^a$ ) to denote the actual states of the motor valves, ( $\tau_1^a, \tau_2^a, \tau_3^a$ ) to time the actual transitions of the motor valves, and ( $P1^a, P2^a, P3^a$ ) to denote the actual states of the pumps. Note that  $Pi^a = w_{ip}$  for all  $i \in \{1, 2, 3\}$ . Furthermore, to relate  $MVi$  to  $w_{im}$ , we consider the following assumption:

(A2) For each  $i \in \{1, 2, 3\}$ , if  $MVi^a \in \{0, 1\}$ , then  $w_{im} \in \{MV_i^a, 2\}$ .

Under (A2), the attacker never bypasses the transition state of the motor valves. Indeed, when the attacker does not respect (A2), the device (valve) goes into a warning state and the attack would be detected immediately (this was verified in the real-world testbed). For example, if  $MV1^a = 1$ , the attacker needs to send first the transition command  $w_{1m} = 2$  before sending the command  $w_{1m} = 0$ .

**Remark 3.** In the absence of actuator attacks, for each  $i \in \{1, 2, 3\}$ , it follows that  $w_{ip} = Pi = Pi^a$  and  $MVi = MV_i^a$ . Furthermore,  $w_{im} = MV_i$  if  $MVi \in \{0, 1\}$  and  $w_{im} = 2$  otherwise. •

### E.1. Control Invariants with Adversary

**Stage 3.** The actual state of the motor valve  $MV3^a$  is governed by the following rules:

$$MV3^a := \begin{cases} 0 & \text{if } MV3^a = 3 \text{ and } \tau_3^a \geq T_3 \\ 1 & \text{if } MV3^a = 2 \text{ and } \tau_3^a \geq T_3 \\ 2 & \text{if } MV3^a = 0 \text{ and } w_{3m} = 2 \\ 3 & \text{if } MV3^a = 1 \text{ and } w_{3m} = 2. \end{cases}$$

Moreover, we reset the value of  $\tau_3^a$  to 0 whenever we switch the value of  $MV3^a$ . In compact form, the discrete behavior of  $MV3^a$  and  $\tau_3^a$  can be modeled by the following constrained difference equation:

$$\begin{aligned} \begin{pmatrix} \tau_3^{a+} \\ MV3^{a+} \end{pmatrix} &= \begin{pmatrix} 0 \\ G_{MV3^a}(MV3^a) \end{pmatrix} \quad (\tau_3^a, MV3^a, w_{3m}) \in D_{MV3^a} \\ D_{MV3^a} &:= \{(\tau_3^a, MV3^a, w_{3m}) : (MV3^a, w_{3m}) \in \{0, 1\} \times \{2\}\} \cup \\ &\quad \{(\tau_3^a, MV3^a, w_{3m}) : (MV3^a, \tau_3^a) \in \{2, 3\} \times \{T_3\}\}, \end{aligned}$$

$$G_{MV3^a}(MV3^a) := \begin{cases} 3 - MV3^a & \text{if } MV3^a = \{2, 3\} \\ 2 + MV3^a & \text{if } MV3^a \in \{0, 1\}. \end{cases}$$

**Stage 2.** Similarly, the discrete behavior of  $MV2^a$  and  $\tau_2^a$  can be modeled by the following constrained difference equation:

$$\begin{pmatrix} \tau_2^{a+} \\ MV2^{a+} \end{pmatrix} = \begin{pmatrix} 0 \\ G_{MV2^a}(MV2^a) \end{pmatrix} \quad (\tau_2^a, MV2^a, w_{2m}) \in D_{MV2^a},$$

$$D_{MV2^a} := \{(\tau_2^a, MV2^a, w_{2m}) : (MV2^a, w_{2m}) \in \{0, 1\} \times \{2\}\} \cup \{(\tau_2^a, MV2^a, w_{2m}) : (MV2^a, \tau_2^a) \in \{2, 3\} \times \{T_2\}\},$$

$$G_{MV2^a}(MV2^a) := \begin{cases} 3 - MV2^a & \text{if } MV2^a \in \{2, 3\} \\ 2 + MV2^a & \text{if } MV2^a \in \{0, 1\}. \end{cases}$$

**Stage 1.** Similarly to Stage 2, the discrete behavior of  $MV1^a$  and  $\tau_1^a$  can be modeled by the following constrained difference equation:

$$\begin{pmatrix} \tau_1^{a+} \\ MV1^{a+} \end{pmatrix} = \begin{pmatrix} 0 \\ G_{MV1^a}(MV1^a) \end{pmatrix} \quad (\tau_1^a, MV1^a, w_{1m}) \in D_{MV1^a},$$

$$D_{MV1^a} := \{(\tau_1^a, MV1^a, w_{1m}) : (MV1^a, w_{1m}) \in \{0, 1\} \times \{2\}\} \cup \{(\tau_1^a, MV1^a, w_{1m}) : (MV1^a, \tau_1^a) \in \{2, 3\} \times \{T_1\}\},$$

$$G_{MV1^a}(MV1^a) := \begin{cases} 3 - MV1^a & \text{if } MV1^a \in \{2, 3\} \\ 2 + MV1^a & \text{if } MV1^a \in \{0, 1\}. \end{cases}$$

### E.2. Physical Invariants with Adversary

Since the attacker can arbitrarily modify the actual state of any pump or motor valve, it follows that (24) is not guaranteed and the plant loses its cascaded interconnection. Indeed, when for example  $MV2 = \text{ON}$ , the stream of the water flowing from Stage 1 to Stage 2 cannot be the same if  $P1 = \text{ON}$  or if  $P1 = \text{OFF}$ . Hence, the dynamics of the water levels  $F_{L3}$  and  $F_{L2}$  in Stages 3 and 2, respectively, will additionally depend on the actual state of pump in the previous stage and the actual state of the motor valve in the next stage. In such a general scenario, the dynamics of the water levels in Stages 3 and 2 can be expressed as follows:

$$\dot{L}_3 = F_{L3}(MV3^a, P2^a), \quad (61)$$

$$\dot{L}_2 = F_{L2}(MV2^a, MV3^a, P2^a, P1^a), \quad (62)$$

where the functions  $F_{L3}$  and  $F_{L2}$  satisfy the following properties:

$$\begin{aligned} F_{L3}(1, 1) &> 0, \quad F_{L3}(2, 1) = F_{L3}(3, 1) > 0, \\ F_{L3}(2, 0) &= F_{L3}(3, 0) = F_{L3}(0, 0) = F_{L3}(0, 1) < 0 \end{aligned} \quad (63)$$

$$\begin{aligned}
F_{L2}(MV2^a, 0, 0, P1^a) &= F_{L2}(MV2^a, 0, 1, P1^a) \\
&= F_{L2}(MV2^a, 2, 0, P1^a) = F_{L2}(MV2^a, 1, 0, P1^a) \geq 0 \\
&\quad \forall (MV2^a, P1^a) \in \{0, 1, 2\} \times \{0, 1\}, \\
F_{L2}(0, MV3^a, P2^a, 0) &= F_{L2}(0, MV3^a, P2^a, 1) \\
&= F_{L2}(2, MV3^a, P2^a, 0) = F_{L2}(1, MV3^a, P2^a, 0) \leq 0 \\
&\quad \forall (MV3^a, P2^a) \in \{0, 1, 2\} \times \{0, 1\}, \\
F_{L2}(MV2^a, MV3^a, 0, 0) &= 0 \quad \forall (MV2^a, MV3^a) \in \{0, 1, 2\}^2, \\
F_{L2}(0, 0, P2^a, P1^a) &= 0 \quad \forall (P2^a, P1^a) \in \{0, 1\}^2, \\
F_{L2}(2, 0, 1, 1) &\geq F_{L2}(2, 2, 1, 1) > 0, \quad F_{L2}(1, 1, 1, 1) > 0, \\
F_{L2}(1, 2, 1, 1) &\geq 0, \quad F_{L2}(2, 1, 1, 1) \leq 0.
\end{aligned} \tag{64}$$

In (64), we are assuming that

$$\begin{aligned}
F_{L2}(MV2^a, 2, P2^a, P1^a) &= F_{L2}(MV2^a, 3, P2^a, P1^a) \\
\text{for all } (MV2^a, P2^a, P1^a) &\in \{0, 1, 2, 3\} \times \{0, 1\} \times \{0, 1\} \\
&\text{and}
\end{aligned}$$

$$\begin{aligned}
F_{L2}(2, MV3^a, P2^a, P1^a) &= F_{L2}(3, MV3^a, P2^a, P1^a) \\
\text{for all } (MV3^a, P2^a, P1^a) &\in \{0, 1, 2, 3\} \times \{0, 1\} \times \{0, 1\}.
\end{aligned}$$

Conditions (63) and (64) are important for our security results in Section 6 to hold. From the numerical values of  $(F_{L3}, F_{L2}, F_{L1})$  in Section B, we can see that these conditions are satisfied for the considered industrial plant.

### E.3. Control+Physical Invariants with Adversary

**Stage 3.** The new augmented state vector for Stage 3 is given by

$$x_3^a := (x_3, \tau_3^a, MV3^a) \in X_3 := X_3 \times [0, T_3] \times \{0, 1, 2, 3\}.$$

Furthermore, the signal sent by the attacker to the third motor valve is  $w_{3m} \in \mathcal{W}_3 := \{0, 1, 2\}$ . Finally, the variable of the second stage affecting the dynamics of the third stage is  $u_3 := P2^a = w_{2p} \in \{0, 1\} := \mathcal{U}_3$ . The dynamical model of Stage 3 under actuator attacks is given by:

$$\mathcal{H}_{3w} : \begin{cases} \dot{x}_3^a = F_{3w}(x_3^a, w_{3m}, u_3) & (x_3^a, w_{3m}, u_3) \in C_{3w} \\ x_3^{a+} = G_{3w}(x_3^a, w_{3m}, u_3) & (x_3^a, w_{3m}, u_3) \in D_{3w}, \end{cases} \tag{65}$$

where  $C_{3w} := \text{cl}((X_3 \times \mathcal{W}_3 \times \mathcal{U}_3) \setminus D_{3w})$ ,

$$\begin{aligned}
D_{3w} &:= \{(x_3^a, w_{3m}, u_3) : x_3 \in D_3\} \cup \\
&\quad \{(x_3^a, w_{3m}, u_3) : (\tau_3^a, MV3^a, w_{3m}) \in D_{MV3^a}\}.
\end{aligned}$$

Furthermore, the jump map  $G_{3w}$  is given by

$$\begin{aligned}
G_{3w}(x_3^a, w_{3m}, u_3) &:= \begin{bmatrix} G_{x_3}(x_3^a, w_{3m}, u_3) \\ G_{\tau_3^a}(x_3^a, w_{3m}, u_3) \\ G'_{MV3^a}(x_3^a, w_{3m}, u_3) \end{bmatrix}, \\
G_{x_3}(x_3^a, w_{3m}, u_3) &:= \begin{cases} G_3(x_3) & \text{if } x_3 \in D_3 \\ x_3 & \text{otherwise,} \end{cases} \\
G'_{MV3^a}(x_3^a, w_{3m}, u_3) &:= \begin{cases} G_{MV3^a}(MV3^a) & \text{if } (\tau_3^a, MV3^a, w_{3m}) \in D_{MV3^a} \\ MV3^a & \text{otherwise,} \end{cases} \\
G_{\tau_3^a}(x_3^a, w_{3m}, u_3) &:= \begin{cases} 0 & \text{if } (\tau_3^a, MV3^a, w_{3m}) \in D_{MV3^a} \\ \tau_3^a & \text{otherwise.} \end{cases}
\end{aligned}$$

Finally, the flow map  $F_{3w}$  is given by

$$F_{3w}(x_3^a, w_{3m}, u_3) := [F_{L3}(MV3^a, w_{2p}) \quad 1 \quad 0 \quad 1 \quad 0]^\top.$$

**Stage 2.** The new augmented state vector is given by

$$x_2^a := (x_2, \tau_2^a, MV2^a) \in X_2 := X_2 \times [0, T_2] \times \{0, 1, 2, 3\}.$$

Furthermore, the signals sent by the attacker to the motor valve and the pump are  $w_2 := (w_{2m}, w_{2p}) \in \mathcal{W}_2 := \{0, 1, 2\} \times \{0, 1\}$ . Finally, the variables of the first and the third stages affecting the dynamics of the second stage are  $u_2 := (L3, P1^a, MV3^a, MV3) \in \mathcal{U}_2$ ,  $\mathcal{U}_2 := \mathbb{R}_{\geq 0} \times \{0, 1\} \times \{0, 1, 2, 3\}^2$ . The dynamical model under actuator attacks is given by:

$$\mathcal{H}_{2w} : \begin{cases} \dot{x}_2^a = F_{2w}(x_2^a, w_2, u_2) & (x_2^a, w_2, u_2) \in C_{2w} \\ x_2^{a+} = G_{2w}(x_2^a, w_2, u_2) & (x_2^a, w_2, u_2) \in D_{2w}, \end{cases} \tag{66}$$

where  $C_{2w} := \text{cl}((X_2 \times \mathcal{W}_2 \times \mathcal{U}_2) \setminus D_{2w})$ ,

$$\begin{aligned}
D_{2w} &:= \{(x_2^a, w_2, u_2) : (x_2, L3, MV3) \in D_2\} \cup \\
&\quad \{(x_2^a, w_{2m}, u_2) : (\tau_2^a, MV2^a, w_{2m}) \in D_{MV2^a}\},
\end{aligned}$$

$$G_{2w}(x_2^a, w_2, u_2) := \begin{bmatrix} G_{x_2}(x_2^a, w_2, u_2) \\ G_{\tau_2^a}(x_2^a, w_2, u_2) \\ G'_{MV2^a}(x_2^a, w_2, u_2) \end{bmatrix},$$

$$\text{where } \begin{aligned} &G_{x_2}(x_2^a, w_2, u_2) := \\ &\begin{cases} G_2(x_2, L3, MV3) & \text{if } (x_2, L3, MV3) \in D_2 \\ x_2 & \text{otherwise,} \end{cases} \end{aligned}$$

$$G'_{MV2^a}(x_2^a, w_2, u_2) := \begin{cases} G_{MV2^a}(MV2^a) & \text{if } (\tau_2^a, MV2^a, w_{2m}) \in D_{MV2^a} \\ MV2^a & \text{otherwise,} \end{cases}$$

$$G_{\tau_2^a}(x_2^a, w_2, u_2) := \begin{cases} 0 & \text{if } (\tau_2^a, MV2^a, w_{2m}) \in D_{MV2^a} \\ \tau_2^a & \text{otherwise,} \end{cases}$$

$$F_{2w}(x_2^a, w_2, u_2) := \begin{bmatrix} F_{L2}(MV2^a, MV3^a, w_{2p}, w_{1p}) \\ 1 \\ 0 \\ 0 \\ 1 \\ 0 \end{bmatrix}.$$

**Stage 3 When P2<sup>a</sup> is Attacked:** In this case, we do not need to extend  $x_3$  nor to consider the signal  $w_{3m}$ . We only consider

$$u_3 := P2^a \in \mathcal{U}_3 := \{0, 1\}$$

affecting the dynamics of Stage 3. Hence, the dynamical model of Stage 3 is given by:

$$\mathcal{H}_3 : \begin{cases} \dot{x}_3 = F_3(x_3, u_3) & (x_3, u_3) \in C_3 \times \mathcal{U}_3 \\ x_3^+ = G_3(x_3) & (x_3, u_3) \in D_3 \times \mathcal{U}_3, \end{cases} \tag{67}$$

where the flow map  $F_3$  is given by

$$F_3(x_3, u) := [F_{L3}(MV3, P2^a) \quad 1 \quad 0]^\top,$$

and  $F_{L3}$  satisfies (63).

**Stage 3 in the Absence of Attacks and When P3 is Hardened:** In this case, the dynamics of L3 can be expressed as follows:

$$\dot{L3} = F_{L3}(MV3, P3).$$

Furthermore, the signal  $P3$  that the controller  $C3$  sends to the third pump is governed by the logic (84). Hence, the local state vector of the third stage is  $\bar{x}_3 := [x_3 \ P3]^\top \in \bar{X}_3$ , where

$$\bar{X}_3 := X_3 \times \{0, 1\}.$$

The hybrid equation modeling Stage 3 is given by:

$$\mathcal{H}_3 : \begin{cases} \dot{\bar{x}}_3 = \bar{F}_3(\bar{x}_3) & \bar{x}_3 \in \bar{C}_3 \\ \bar{x}_3^+ = \bar{G}_3(\bar{x}_3) & \bar{x}_3 \in \bar{D}_3, \end{cases} \quad (68)$$

where

$$\begin{aligned} \bar{C}_3 &:= \text{cl}(\bar{X}_3 \setminus \bar{D}_3), \\ \bar{D}_3 &:= (D_3 \times \{0, 1\}) \cup \{\bar{x}_3 : (L3, P3) \in D_{P3}\}. \end{aligned}$$

Furthermore, the jump map  $\bar{G}_3$  is given by:

$$\bar{G}_3(\bar{x}_3) := [G_{x_3}(\bar{x}_3) \ G'_{P3}(\bar{x}_3)]^\top,$$

where

$$\begin{aligned} G_{x_3}(\bar{x}_3) &:= \begin{cases} G_3(x_3) & \text{if } x_3 \in D_3 \\ x_3 & \text{otherwise,} \end{cases} \\ G'_{P3}(\bar{x}_3) &:= \begin{cases} 1 - P3 & \text{if } (L3, P3) \in D_{P3} \\ P3 & \text{otherwise,} \end{cases} \end{aligned}$$

and

$$\bar{F}_3(\bar{x}_3) := [F_{L3}(MV3, P3) \ 1 \ 0 \ 0]^\top,$$

with

$$\begin{aligned} F_{L3}(2, 0) = F_{L3}(3, 0) &> 0, \quad F_{L3}(2, 1) = F_{L3}(3, 1) < 0, \\ F_{L3}(1, P3) &> 0 \quad \forall P3 \in \{0, 1\}, \quad F_{L3}(0, 1) < 0, \\ F_{L3}(0, 0) &= 0. \end{aligned} \quad (69)$$

**Stage 3 Under Actuator Attacks When  $P3$  is Hardened:** In the following, we extend the dynamical model in (65) when  $P3$  is not Always Open. In such a general scenario, the flow dynamics (61) in the third stage depends also on  $P3$  and can be expressed as follows:

$$\dot{L3} = F_{L3}(MV3^a, P3^a, P2^a). \quad (70)$$

Furthermore, we still assume that the attacker affects only the actual state of the third motor valve. Hence, the new state vector is given by

$$\bar{x}_3^a := (\bar{x}_3, \tau_3^a, MV3^a) \in \bar{X}_3^a := \bar{X}_3 \times [0, T_3] \times \{0, 1, 2, 3\}.$$

Furthermore, the signal sent by the attacker to the third motor valve is

$$w_{3m} \in \mathcal{W}_3 := \{0, 1, 2\}.$$

Finally, the variable of the second stage affecting the dynamics of the third stage is

$$u_3 := P2^a \in \mathcal{U}_3 := \{0, 1\}.$$

The dynamical model of Stage 3 under actuator attacks is given by:

$$\mathcal{H}_{3w} : \begin{cases} \dot{\bar{x}}_3^a = \bar{F}_{3w}(\bar{x}_3^a, w_{3m}, u_3) & (\bar{x}_3^a, w_{3m}, u_3) \in \bar{C}_{3w} \\ \bar{x}_3^{a+} = \bar{G}_{3w}(\bar{x}_3^a, w_{3m}, u_3) & (\bar{x}_3^a, w_{3m}, u_3) \in \bar{D}_{3w}, \end{cases}$$

where the sets  $\bar{C}_{3w}$  and  $\bar{D}_{3w}$  are given by

$$\begin{aligned} \bar{C}_{3w} &:= \text{cl}((\bar{X}_3^a \times \mathcal{W}_3 \times \mathcal{U}_3) \setminus \bar{D}_{3w}), \\ \bar{D}_{3w} &:= \{(\bar{x}_3^a, w_{3m}, u_3) : \bar{x}_3 \in \bar{D}_3\} \cup \\ &\quad \{(\bar{x}_3^a, w_{3m}, u_3) : (\tau_3^a, MV3^a, w_{3m}) \in D_{MV3^a}\}, \end{aligned}$$

$$\bar{G}_{3w}(\bar{x}_3^a, w_{3m}, u_3) := \begin{bmatrix} G_{\bar{x}_3}(\bar{x}_3^a, w_{3m}, u_3) \\ G'_{\tau_3^a}(\bar{x}_3^a, w_{3m}, u_3) \\ G'_{MV3^a}(\bar{x}_3^a, w_{3m}, u_3) \end{bmatrix},$$

where

$$G_{\bar{x}_3}(\bar{x}_3^a, w_{3m}, u_3) := \begin{cases} \bar{G}_3(\bar{x}_3) & \text{if } \bar{x}_3 \in \bar{D}_3 \\ \bar{x}_3 & \text{otherwise,} \end{cases}$$

$$G'_{MV3^a}(\bar{x}_3^a, w_{3m}, u_3) := \begin{cases} G_{MV3^a}(MV3^a) & \text{if } (\tau_3^a, MV3^a, w_{3m}) \in D_{MV3^a} \\ MV3^a & \text{otherwise,} \end{cases}$$

$$G'_{\tau_3^a}(\bar{x}_3^a, w_{3m}, u_3) := \begin{cases} 0 & \text{if } (\tau_3^a, MV3^a, w_{3m}) \in D_{MV3^a} \\ \tau_3^a & \text{otherwise.} \end{cases}$$

Finally, the flow map  $\bar{F}_{3w}$  is given by

$$\bar{F}_{3w}(\bar{x}_3^a, w_{3m}, u) := [F_{L3}(MV3^a, P3^a, P2^a) \ 1 \ 0 \ 0 \ 1 \ 0]^\top,$$

where

$$\begin{aligned} F_{L3}(0, P3^a, 0) &= F_{L3}(0, P3^a, 1) = F_{L3}(2, P3^a, 0) = \\ F_{L3}(3, P3^a, 0) &= F_{L3}(1, P3^a, 0) \leq 0, \end{aligned}$$

$$F_{L3}(MV3^a, 0, P2^a) \geq 0,$$

$$F_{L3}(MV3^a, 0, 0) = 0 \quad \forall MV3^a \in \{0, 1, 2, 3\},$$

$$F_{L3}(0, 0, P2^a) = 0 \quad \forall P2^a \in \{0, 1\},$$

$$F_{L3}(MV3^a, 1, 1) > 0 \quad \forall MV3^a \in \{2, 3\},$$

$$F_{L3}(1, 1, 1) > 0.$$

**Stage 3 When  $P2^a$  is Attacked and When  $P3$  is Hardened:** In this case, we include  $u_3 := P2^a = w_{2m} \in \mathcal{U}_3 := \{0, 1\}$  as an external signal affecting Stage 3. Hence, the dynamical model of Stage 3 when  $P2^a$  is attacked and  $P3$  is not Always Open is given by:

$$\mathcal{H}_3 : \begin{cases} \dot{\bar{x}}_3 = \bar{F}_3(\bar{x}_3, u_3) & (\bar{x}_3, u_3) \in \bar{C}_3 \times \mathcal{U}_3 \\ \bar{x}_3^+ = \bar{G}_3(\bar{x}_3) & (\bar{x}_3, u_3) \in \bar{D}_3 \times \mathcal{U}_3, \end{cases} \quad (71)$$

where

$$\bar{F}_3(\bar{x}_3, u_3) := [F_{L3}(MV3, P3, P2^a) \ 1 \ 0 \ 0]^\top$$

with

$$\begin{aligned} F_{L3}(0, P3^a, 0) &= F_{L3}(0, P3^a, 1) = \\ F_{L3}(2, P3^a, 0) &= F_{L3}(3, P3^a, 0) = \\ &= F_{L3}(1, P3^a, 0) \leq 0 \quad \forall P3^a \in \{0, 1\}, \\ F_{L3}(MV3^a, 0, P2^a) &\geq 0 \quad \forall (MV3^a, P2^a) \in \{0, 1, 2, 3\} \times \{0, 1\}, \\ F_{L3}(MV3^a, 0, 0) &= 0 \quad \forall MV3^a \in \{0, 1, 2, 3\}, \\ F_{L3}(0, 0, P2^a) &= 0 \quad \forall P2^a \in \{0, 1\}, \\ F_{L3}(MV3^a, 1, 1) &> 0 \quad \forall MV3^a \in \{2, 3\}, \\ F_{L3}(1, 1, 1) &> 0. \end{aligned} \quad (72)$$

**Stage 2 When  $P2^a$  is Attacked:** In this case, we do not need to extend the state vector  $x_2$ . Furthermore, since the actuator attack affects  $P2^a$  only; it follows that

$$w_2 := w_{2p} \in \mathcal{W}_2 := \{0, 1\}.$$

Finally, the variables from Stages 1 and 3 that affect the dynamics of Stage 2 are

$$\begin{aligned} u_2 &:= (L3, MV3) \in \mathcal{U}_2 \\ \mathcal{U}_2 &:= [L3_{\min}, L3_{\max} + \delta] \times \{0, 1, 2, 3\}. \end{aligned}$$

The dynamical model under actuator attacks is given by

$$\mathcal{H}_{2w} : \begin{cases} \dot{x}_2 = F_{2w}(x_2, w_2, u_2) & (x_2, w_2, u_2) \in C_{2w} \\ x_2^+ = G_2(x_2, u_2) & (x_2, w_2, u_2) \in D_{2w}, \end{cases} \quad (73)$$

where the sets  $C_{2w}$  and  $D_{2w}$  are given by

$$\begin{aligned} C_{2w} &:= \text{cl}((X_2 \times \mathcal{W}_2 \times \mathcal{U}_2) \setminus D_{2w}), \\ D_{2w} &:= \{(x_2, w_2, u_2) \in X_2 \times \mathcal{W}_2 \times \mathcal{U}_2 : (x_2, u_2) \in D_2\}. \end{aligned}$$

Finally, the flow map  $F_{2w}$  is given by

$$F_{2w}(x_2, w_2, u_2) := \begin{bmatrix} F_{L2}(MV2, MV3, P2^a = w_{2p}) \\ 1 \\ 0 \\ 0 \end{bmatrix}.$$

Using (64), and (24) when  $i = 1$ , we conclude that

$$\begin{aligned} F_{L2}(MV2^a, 0, 0) &= F_{L2}(MV2^a, 0, 1) = F_{L2}(MV2^a, 2, 0) \\ &= F_{L2}(MV2^a, 1, 0) \geq 0 \quad \forall MV2^a \in \{0, 1, 2\}, \\ F_{L2}(0, 0, P2^a) &= 0 \quad \forall P2^a \in \{0, 1\}, \\ F_{L2}(2, 0, 1) &\geq F_{L2}(2, 2, 1) > 0, \quad F_{L2}(1, 1, 1) > 0, \\ F_{L2}(1, 2, 1) &\geq 0, \quad F_{L2}(2, 1, 1) \leq 0. \end{aligned} \quad (74)$$

**Remark 4.** Note that, in (74), we are assuming that

$$\begin{aligned} F_{L2}(MV2^a, 2, P2^a) &= F_{L2}(MV2^a, 3, P2^a) \\ \text{for all } (MV2^a, P2^a) &\in \{0, 1, 2, 3\} \times \{0, 1\} \text{ and} \\ F_{L2}(2, MV3^a, P2^a) &= F_{L2}(3, MV3^a, P2^a) \\ \text{for all } (MV3^a, P2^a) &\in \{0, 1, 2, 3\} \times \{0, 1\}. \end{aligned} \quad \bullet$$

**Remark 5.** The systems in (67) and (73) are particular cases of (65) and (66), respectively, when the motor valves are not attacked; namely, when  $MV_i = MV_i^a$  and  $\tau_i = \tau_i^a$  for all  $i \in \{2, 3\}$ , and when (24) is satisfied for the first stage; namely, for  $i = 1$ .  $\bullet$

**Stage 2 When  $P1^a$  is Attacked:** In this case, we don't need to extend the state vector  $x_2$ . Furthermore, the variables of the first and the third stages affecting the dynamics of the second stage are

$$\begin{aligned} u_2 &:= (L3, MV3, P1^a) \in \mathcal{U}_2 \\ \mathcal{U}_2 &:= [L3_{\min}, L3_{\max} + \delta] \times \{0, 1, 2, 3\} \times \{0, 1\}. \end{aligned}$$

The dynamical model is given by:

$$\mathcal{H}_2 : \begin{cases} \dot{x}_2 = F_2(x_2, u_2) & (x_2, u_2) \in C_2 \times \{0, 1\} \\ x_2^+ = G_2(x_2, L3, MV3) & (x_2, u_2) \in D_2 \times \{0, 1\}, \end{cases} \quad (75)$$

where the flow map  $F_2$  is given by

$$F_2(x_2, u_2) := \begin{bmatrix} F_{L2}(MV2, P2, P1^a) \\ 1 \\ 0 \\ 0 \end{bmatrix}, \quad (76)$$

with

$$\begin{aligned} F_{L2}(MV2^a, 0, P1^a) &\geq 0 \quad \forall (MV2^a, P1^a) \in \{0, 1, 2, 3\} \times \{0, 1\}, \\ F_{L2}(0, P2^a, 0) &= F_{L2}(0, P2^a, 1) = F_{L2}(2, P2^a, 0) = \\ &= F_{L2}(3, P2^a, 0) = \\ &= F_{L2}(1, P2^a, 0) \leq 0 \quad \forall P2^a \in \{0, 1\}, \\ F_{L2}(1, 1, 1) &> 0, \\ F_{L2}(MV2^a, 0, 0) &= 0 \quad \forall MV2^a \in \{0, 1, 2, 3\}, \\ F_{L2}(0, 0, P1^a) &= 0 \quad \forall P1^a \in \{0, 1\}, \\ F_{L2}(MV2^a, 1, 1) &\leq 0 \quad \forall MV2^a \in \{2, 3\}. \end{aligned} \quad (77)$$

**Stage 2 in the Absence of Attacks When the Control Logic of  $P2$  is modified:** According to (87) and compared to (16), the pump  $P2$  closes when  $L2 \leq L2_{\min} - \delta$ . Hence, we extend the model (26) as follows:

$$\mathcal{H}_2 : \begin{cases} \dot{x}_2 = F_2(x_2) & (x_2, u_2) \in \tilde{C}_2 \\ x_2^+ = \tilde{G}_2(x_2, u_2) & (x_2, u_2) \in \tilde{D}_2, \end{cases} \quad (78)$$

where

$$\begin{aligned} u_2 &:= (L3, MV3) \in \mathcal{U}_2 \\ \mathcal{U}_2 &:= [L3_{\min}, L3_{\max} + \delta] \times \{0, 1, 2, 3\}, \\ \tilde{C}_2 &:= \text{cl}((X_2 \times \mathcal{U}_2) \setminus \tilde{D}_2), \quad \tilde{D}_2 := \tilde{D}_{21} \cup \tilde{D}_{22}, \end{aligned}$$

$$\begin{aligned} \tilde{D}_{21} &:= \{(x_2, u_2) : L2 \geq L2_{\max}, MV2 = 1\} \cup \\ &\quad \{(x_2, u_2) : \tau_2 \geq T_2, MV2 \in \{2, 3\}\} \cup \\ &\quad \{(x_2, u_2) : L2 \leq L2_{\min}, MV2 = 0\}, \\ \tilde{D}_{22} &:= \{(x_2, u_2) : L3 \leq L3_{\min}, P2 = 0, L2 \geq L2_{\min}\} \cup \\ &\quad \{(x_2, u_2) : MV3 \in \{1, 2\}, P2 = 0, L2 \geq L2_{\min}\} \cup \\ &\quad \{(x_2, u_2) : L3 \leq L3_{\min}, P2 = 0, L2 \geq L2_{\min}\} \cup \\ &\quad \{(x_2, u_2) : MV3 \in \{1, 2\}, P2 = 0, L2 \geq L2_{\min}\} \cup \\ &\quad \{(x_2, u_2) : L3 \geq L3_{\max}, P2 = 1\} \cup \\ &\quad \{(x_2, u_2) : MV3 \in \{0, 3\}, P2 = 1\} \cup \\ &\quad \{(x_2, u_2) : L2 \leq L2_{\min}, P2 = 1\}. \end{aligned}$$

Furthermore, the jump map  $\tilde{G}_2$  is given by:

$$\tilde{G}_2(x_2, u_2) := \begin{bmatrix} L2 \\ G_{22}(x_2, u_2) \\ G_{23}(x_2, u_2) \\ \tilde{G}_{24}(x_2, u_2) \end{bmatrix},$$

where

$$G_{22}(x_2, u_2) := \begin{cases} 0 & \text{if } (x_2, u_2) \in \tilde{D}_{21} \\ \tau_2 & \text{otherwise,} \end{cases}$$

$$G_{23}(x_2, u_2) := \begin{cases} G_{MV2}(MV2) & \text{if } (x_2, u_2) \in \tilde{D}_{21} \\ MV2 & \text{otherwise,} \end{cases}$$

and

$$\tilde{G}_{24}(x_2, u_2) := \begin{cases} 1 - P2 & \text{if } (x_2, u_2) \in \tilde{D}_{22} \\ P2 & \text{otherwise.} \end{cases}$$

**Stage 2 When  $P1^a$  is Attacked and the Control Logic of  $P2$  is Modified:** The dynamical model is given by:

$$\mathcal{H}_2 : \begin{cases} \dot{x}_2 = F_2(x_2, u_2) & (x_2, u_2) \in \tilde{C}_2 \times \{0, 1\} \\ x_2^+ = \tilde{G}_2(x_2, L3, MV3) & (x_2, u_2) \in \tilde{D}_2 \times \{0, 1\}, \end{cases} \quad (79)$$

where

$$u_2 := (\text{L3}, \text{MV3}, \text{P1}^a) \in \mathcal{U}_2,$$

$$\mathcal{U}_2 := [\text{L3min}, \text{L3max} + \delta] \times \{0, 1, 2, 3\} \times \{0, 1\},$$

and  $F_2$  as in (76).

## F. Security Proofs Under Actuation Attacks

**Safety With Attacks:** As our first contribution, we now adapt our previous results to reason about safety under attacks. First, we consider a hybrid system under general attacks

$$\mathcal{H}_u : \begin{cases} \dot{x} = F(x, u) & (x, u) \in C \\ x^+ = G(x, u) & (x, u) \in D. \end{cases} \quad (80)$$

Where the attack  $u$  can affect the physical states, as well as the discrete software logic. To analyze safety in the presence of attacks, we introduce a new concept we call uniform safety.

**Definition 8 (Uniform Safety).** System  $\mathcal{H}_u$  in (80) is said to be safe with respect to  $(X_o, X_u)$  uniformly in  $u \in \mathcal{U}$  iff, for each solution pair  $(x, u)$  to  $\mathcal{H}_u$  such that  $x(0, 0) \in X_o$ , the solution  $x$  never reaches the set  $X_u$ .

In the presence of attacks, the variables  $(\text{Pi}, \text{MVi})$  do not necessarily correspond to the actual states of the  $i$ -th motor valve and the  $i$ -th pump, respectively. For this reason, we introduce the extra variables  $(\text{MV1}^a, \text{MV2}^a, \text{MV3}^a)$  to denote the actual states of the motor valves,  $(\tau_1^a, \tau_2^a, \tau_3^a)$  to time the actual transitions of the motor valves, and  $(\text{P1}^a, \text{P2}^a, \text{P3}^a)$  to denote the actual states of the pumps.

As we show in the experimental results, the system with the original PLC programs is unsafe in the presence of attacks. Indeed, due to the constant demand of water by Stage 3,  $\text{L2}$  becomes less than  $\text{L2min} - \delta$  if  $\text{P1}^a = 0$  is maintained by the attacker. Similarly,  $\text{L3}$  becomes less than  $\text{L3min}$  if  $\text{P2}^a = 0$  is maintained. Therefore the original system is unsafe to attacks that can compromise either the first or the second pump. However, as we will show in this section, if we change the control logic of PLCs, the system can be made safe against arbitrary attacks (as long as they compromise only one control signal).

**Claim 3.** When the attacker forces  $\text{P1}^a = 0$  (closes the pump in Stage 1) and/or forces  $\text{P2}^a = 0$  (closes the pump in Stage 2), the plant becomes unsafe. ■

We prove Claim 3 by finding a counterexample. Indeed, due to the constant demand of water by Stage 3,  $\text{L2}$  becomes less than  $\text{L2min} - \delta$  if  $\text{P1}^a = 0$  is maintained by the attacker. Similarly,  $\text{L3}$  becomes less than  $\text{L3min}$  if  $\text{P2}^a = 0$  is maintained.

**Claim 4.** The plant remains safe if the attacker forces  $\text{P2}^a = 1$  and/or  $\text{P1}^a = 1$ . ■

To prove Claim 4, we will formally show that the plant remains safe when the attacker enforces  $\text{P2}^a = 1$ . Intuitively, in this case, the flow of water is governed by the motor valves and the behavior of the plant is not very different from its behavior in the absence of attacks.

The only change concerns the dynamics of  $(\text{L1}, \text{L2}, \text{L3})$ , which as we shall show, does not compromise the safety of the plant. We model Stages 3 and 2 when  $\text{P2}^a$  is attacked; see system  $\mathcal{H}_3$  and system  $\mathcal{H}_2$  Theorems 10 and 11 prove Claim 2 by showing safety of the plant when the attacker forces  $\text{P2}^a = 1$ .

**Theorem 10.** Consider the hybrid system  $\mathcal{H}_3$  with  $u_3 = \text{P2}^a = 1$ . Consider the initial set  $X_{o3}$  in (29) and an unsafe set  $X_{u3} \subset X_3 \setminus X_{s3}$  with  $X_{s3}$  introduced in (35). Assume that there exists  $\sigma > 0$  such that

$$4T_3(F_{\text{L3}}(3, 1) + \sigma) \leq \delta. \quad (81)$$

Then, the hybrid system  $\mathcal{H}_3$  with  $u_3 = \text{P2}^a = 1$  is safe with respect to  $(X_{o3}, X_{u3})$ , and admits a barrier function certificate given by

$$B(x_3) := (\text{L3} - \text{L3min})(\text{L3} - \text{L3max} - f(\tau_3, \text{MV3})),$$

where  $f(\tau_3, \text{MV3}) := (F_{\text{L3}}(3, 1) + \sigma)[\tau_3 + T_3 * w_f(\text{MV3})]$ , where  $w_f(2) := 0$ ,  $w_f(1) := 1$ ,  $w_f(3) := 2$ ,  $w_f(0) := 3$ . ■

*Proof.* It can be derived by following the exact same steps as in the proof of Theorem 6 while noting, under (63), that

$$F_{\text{L3}}(1, 1) > 0, F_{\text{L3}}(2, 1) = F_{\text{L3}}(3, 1) > 0, F_{\text{L3}}(0, 1) < 0. \quad \blacksquare$$

**Theorem 11.** Consider the hybrid system  $\mathcal{H}_2$ . Consider the initial set  $X_{o2}$  in (30) and the unsafe set  $X_{u2} \subset X_2 \setminus X_{s2}$  with  $X_{s2}$  introduced in (34). Assume that there exist  $\sigma_h > 0$  and  $\sigma_g > 0$  such that

$$4T_2(F_{\text{L2}}(2, 0, 0) + \sigma_h) \leq \delta, \quad (82)$$

$$4T_2(F_{\text{L2}}(2, 1, 1) + \sigma_g) \leq \delta. \quad (83)$$

Then, the hybrid system  $\mathcal{H}_2$  is safe with respect to  $(X_{o2}, X_{u2})$  uniformly in  $(u_2, w_2) \in \mathcal{U}_2 \times \mathcal{W}_2$ , and admits a barrier function certificate given by

$$B(x_2) := (\text{L2} - \text{L2min} + g(\tau_2, \text{MV2})) \times (\text{L2} - \text{L2max} - h(\tau_2, \text{MV2})),$$

where  $g(\tau_2, \text{MV2}) := (-F_{\text{L2}}(2, 1, 1) + \sigma_g) * [\tau_2 + T_2 * w_g(\text{MV2})]$ ,  $w_g(3) := 0$ ,  $w_g(0) := 1$ ,  $w_g(2) := 2$ ,  $w_g(1) := 3$ ,  $h(\tau_2, \text{MV2}) := (F_{\text{L2}}(2, 0, 0) + \sigma_h) * [\tau_2 + T_2 * w_h(\text{MV2})]$ ,  $w_h(2) := 0$ ,  $w_h(1) := 1$ ,  $w_h(3) := 2$ ,  $w_h(0) := 3$ . ■

*Proof.* We note that the system in (73) can be expressed as follows:

$$\mathcal{H}_{2w} : \begin{cases} \dot{x}_2 = F_{2w}(x_2, w_2, u_2) \\ (x_2, u_2, w_2) \in C_2 \times \mathcal{W}_2 \\ x_2^+ = G_2(x_2, u_2) \\ (x_2, u_2, w_2) \in D_2 \times \mathcal{W}_2. \end{cases}$$

The rest of the proof follows exactly using the same steps as in the proof of Theorem 7 while noting, under (74), that

$$F_{\text{L2}}(\text{MV2}, \text{MV3}, \text{P2}^a) \geq F_{\text{L2}}(2, 1, 1)$$

for all  $(\text{MV2}, \text{MV3}, \text{P2}^a) \in \{1, 2\} \times \{0, 1, 2\} \times \{0, 1\}$ , and

$$F_{\text{L2}}(\text{MV2}, \text{MV3}, \text{P2}^a) \leq F_{\text{L2}}(2, 0, 0)$$

for all  $(\text{MV2}, \text{MV3}, \text{P2}^a) \in \{0, 2\} \times \{0, 1, 2\} \times \{0, 1\}$ . ■



## F.1. Changing the Control Logic of PLCs to Make the System More Secure

In this section we harden the system to make it more resilient to attacks. In particular, we first change the control logic of the PLC controlling stage 3 (C3) so that P3 is not always 1. As a result, we include P3 as a control parameter governed by the following logic:

$$P3 := \begin{cases} 0 & \text{if } (L3 \leq L3_o, P3 = 1) \\ 1 & \text{if } (L3 \geq L3_o, P3 = 0), \end{cases} \quad (84)$$

where  $L3_o > 0$  is a lower bound on the water level  $L3$  in Stage 3, it aims to avoid the dry-runs (operates without liquid) of the pump P3. Hence, the behavior of P3 can be modeled by the following constrained difference equation:

$$P3^+ = G_{P3}(P3) \quad (L3, P3) \in D_{P3},$$

where  $D_{P3} := \{(L3, P3) : L3 \leq L3_o, P3 = 1\} \cup \{(L3, P3) : L3 \geq L3_o, P3 = 0\}$  and  $G_{P3}(P3) := 1 - P3$ .

Using the logic (84) with  $L3_o = L3_{\min}$ , we are able to show the following claim.

**Claim 5.** When P3 is governed by (84) with  $L3_o = L3_{\min}$ , the plant remains safe under any arbitrary time series of possible attacks affecting  $P2^a$ . •

To show Claim 5, we use Theorem 11 to conclude that it is enough to show safety of Stage 3 uniformly in  $P2^a \in \{0, 1\}$  when P3 is governed by (84). To simplify the analysis, we model Stage 3 when only  $P2^a$  is attacked and P3 is controlled by our modified control logic (it is hardened); see system  $\mathcal{H}_3$ .

**Theorem 12.** Consider the hybrid system  $\mathcal{H}_3$ . Consider the initial set  $\bar{X}_{o3} := X_{o3} \times \{0, 1\}$  and an unsafe set  $X_{u3} \subset \bar{X}_3 \setminus \bar{X}_{s3}$  with  $\bar{X}_{s3} = X_{s3} \times \{0, 1\}$ . Assume that there exists  $\sigma > 0$  such that

$$4T_3(F_{L3}(3, 0, 1) + \sigma) \leq \delta. \quad (85)$$

Then, the hybrid system  $\mathcal{H}_3$  is safe with respect to  $(\bar{X}_{o3}, \bar{X}_{u3})$  uniformly in  $u_3 = w_{2m} = P2^a \in \mathcal{U}_3$ , and admits a barrier function certificate given by

$$B(\bar{x}_3) := (L3 - L3_{\min})(L3 - L3_{\max} - P3 * f(\tau_3, MV3)),$$

where  $f(\tau_3, MV3) := (F_{L3}(3, 0, 1) + \sigma)[\tau_3 + T_3 * w_f(MV3)]$ , and  $w_f(2) := 0$ ,  $w_f(1) := 1$ ,  $w_f(3) := 2$ ,  $w_f(0) := 3$ . □

*Proof.*

Note that

$$\bar{D}_3 := \bar{D}_{31} \cup (D_3 \times \{0, 1\}),$$

where

$$\bar{D}_{31} := \{\bar{x}_3 \in \bar{X}_3 : L3 \leq L3_{\min}, P3 = 1\} \cup \{\bar{x}_3 \in \bar{X}_3 : L3 \geq L3_{\min}, P3 = 0\}.$$

Furthermore, we introduce the set

$$\begin{aligned} \bar{C}_3 &:= \text{cl}(\bar{X}_3 \setminus (\bar{D}_{31} \cup (D_3 \times \{0, 1\}))), \\ &= \text{cl}(\bar{X}_3 \setminus \bar{D}_{31}) \cap \text{cl}(\bar{X}_3 \setminus (D_3 \times \{0, 1\})) \\ &= \text{cl}(\bar{X}_3 \setminus \bar{D}_{31}) \cap (\text{cl}(X_3 \setminus D_3) \times \{0, 1\}) \\ &= \text{cl}(\bar{X}_3 \setminus \bar{D}_{31}) \cap (C_3 \times \{0, 1\}). \end{aligned}$$

Finally, we let

$$\begin{aligned} \bar{C}_{31} &:= \text{cl}(\bar{X}_3 \setminus \bar{D}_{31}) \\ &= \{\bar{x}_3 \in \bar{X}_3 : L3 \geq L3_{\min}, P3 = 1\} \cup \\ &\quad \{\bar{x}_3 \in \bar{X}_3 : L3 \leq L3_{\min}, P3 = 0\}, \end{aligned}$$

to conclude that

$$\bar{C}_3 = \bar{C}_{31} \cap (C_3 \times \{0, 1\}).$$

Furthermore, note that

$$f(\tau_3, MV3) \in [0, (F_{L3}(3, 0, 1) + \sigma) * 4 * T_3]$$

for all  $(\tau_3, MV3) \in [0, T_3] \times \{0, 1, 2, 3\}$ . Hence, for  $\sigma > 0$  and  $\delta > 0$  such that (85) holds, we conclude that

$$f(\tau_3, MV3) \in [0, \delta] \quad \forall (\tau_3, MV3) \in [0, T_3] \times \{0, 1, 2, 3\}.$$

Thus, (4) is satisfied. Next, we let the set

$$K_e := \{\bar{x}_3 \in \bar{X}_3 : L3 \in [L3_{\min}, L3_{\max} + P3 * f(\tau_3, MV3)]\}.$$

To complete the proof, we use Theorem 1 and we start verifying the jump conditions (11) and (12). Note that the set  $K_e \cap \bar{D}_3$  satisfies

$$K_e \cap \bar{D}_3 = A_1 \cup A_2 \cup A_3 \cup A_4,$$

where

$$\begin{aligned} A_1 &:= \{\bar{x}_3 \in \bar{X}_3 : L3 \in [L3_{\max}, L3_{\max} + P3 * f(\tau_3, 1)], \\ &\quad MV3 = \{1\}, P3 = 1\}, \\ A_2 &:= \{\bar{x}_3 \in \bar{X}_3 : L3 \in [L3_{\min}, L3_{\max} + P3 * f(T_3, MV3)], \\ &\quad \tau_3 = T_3, MV3 \in \{2, 3\}, P3 = 1\}, \\ A_3 &:= \{\bar{x}_3 \in \bar{X}_3 : L3 = L3_{\min}, P3 = 1\}, \\ A_4 &:= \{\bar{x}_3 \in \bar{X}_3 \setminus (A_1 \cup A_2 \cup A_3) : \\ &\quad L3 \in [L3_{\min}, L3_{\max} + P3 * f(T_3, MV3)], P3 = 0\}. \end{aligned}$$

Note that, for each  $\bar{x}_3 \in A_1$ , we have

$$\bar{G}_3(\bar{x}_3) = [L3 \ 0 \ 3 \ P3]^T.$$

Hence,

$$\begin{aligned} B(\bar{G}_3(\bar{x}_3)) &= \\ (L3 - L3_{\min})(L3 - L3_{\max} - P3 * f(0, 3)) &\leq 0. \end{aligned}$$

The latter inequality is true since

$$f(\tau_3, 1) \leq f(0, 3) \quad \forall \tau_3 \in [0, T_3].$$

Similarly, for each  $\bar{x}_3 \in A_2$ , we have

$$\bar{G}_3(\bar{x}_3) = [L3 \ 0 \ \alpha(MV3) \ P3]^T,$$

where  $\alpha(3) := 0$  and  $\alpha(2) := 1$ . Hence,

$$\begin{aligned} B(\bar{G}_3(\bar{x}_3)) &= \\ (L3 - L3_{\min})(L3 - L3_{\max} - P3 * f(0, \alpha(MV3))) &\leq 0. \end{aligned}$$

The latter inequality is true since

$$f(\tau_3, 3) \leq f(0, 0) \quad \forall \tau_3 \in [0, T_3],$$

and

$$f(\tau_3, 2) \leq f(0, 1) \quad \forall \tau_3 \in [0, T_3].$$

Next, for each  $\bar{x}_3 \in A_3$ , we have

$$\bar{G}_3(\bar{x}_3) = [L3_{\min} \ \star \ \star \ 0]^T.$$

Hence,  $B(\bar{G}_3(\bar{x}_3)) = 0$ . Finally, for each  $\bar{x}_3 \in A_4$ , we have

$$\bar{G}_3(\bar{x}_3) = [\text{L3} \quad \star \quad \star \quad 1]^\top,$$

$$B(\bar{x}_3) = (\text{L3} - \text{L3min})(\text{L3} - \text{L3max}) \leq 0.$$

Hence,

$$B(\bar{G}_3(\bar{x}_3)) = (\text{L3} - \text{L3min})(\text{L3} - \text{L3max} - f(\tau_3, \text{MV3})) \leq 0.$$

We conclude that (11) is satisfied. Moreover, to show (12), we notice that  $\bar{C}_3 \cup \bar{D}_3 = \bar{X}_3$  and  $\bar{G}_3(\bar{x}_3) \in \bar{X}_3$  for all  $\bar{x}_3 \in \bar{D}_3$ .

Next, to verify (10), we start noting that the set  $U(\partial K_e) \setminus K_e$  satisfies

$$U(\partial K_e) \setminus K_e = B_1 \cup B_2,$$

where, for some  $\epsilon > 0$ ,

$$B_1 := \{\bar{x}_3 \in \bar{X}_3 : \text{L3} \in (\text{L3min} - \epsilon, \text{L3min})\},$$

and

$$B_2 := \{\bar{x}_3 \in \bar{X}_3 : \text{L3} \in (\text{L3max} + \text{P3} * f(\tau_3, \text{MV3}))[1, \epsilon)\}.$$

Furthermore, the set  $\bar{C}_3$  can be explicitly expressed as

$$\bar{C}_3 = (C_3 \times \{0, 1\}) \cap \bar{C}_{31} = \left( \bigcap_{i=1}^3 \bar{D}_{3i} \right) \cap \bar{C}_{31},$$

where

$$\bar{D}_{31} := \{\bar{x}_3 \in \bar{X}_3 : \text{L3} \geq \text{L3min} \cup \text{MV3} \in \{1, 2, 3\}\}$$

$$\bar{D}_{32} := \{\bar{x}_3 \in \bar{X}_3 : \tau_3 \leq T_3 \cup \text{MV3} \in \{0, 1\}\}$$

$$\bar{D}_{33} := \{\bar{x}_3 \in \bar{X}_3 : \text{L3} \leq \text{L3max} \cup \text{MV3} \in \{0, 2, 3\}\},$$

This is equivalent to

$$\bar{C}_3 := \left( \bigcup_{i=1}^4 C_{3i} \right) \cap \bar{C}_{31}, \quad (86)$$

where

$$C_{31} := \{\bar{x}_3 \in \bar{X}_3 : \text{L3min} \leq \text{L3} \leq \text{L3max}\}$$

$$C_{32} := \{\bar{x}_3 \in \bar{X}_3 : \text{L3} \leq \text{L3max}, \text{MV3} \in \{1, 2, 3\}\}$$

$$C_{33} := \{\bar{x}_3 \in \bar{X}_3 : \text{L3} \geq \text{L3min}, \text{MV3} \in \{0, 2, 3\}\}$$

$$C_{34} := \{\bar{x}_3 \in \bar{X}_3 : \text{MV3} \in \{2, 3\}\}.$$

We can also show that

$$\bar{C}_3 := C_{31}^a \cup C_{31}^b \cup C_{32}^a \cup C_{32}^b \cup C_{33}^a \cup C_{33}^b \cup C_{34}^a \cup C_{34}^b,$$

where

$$C_{31}^a := \{\bar{x}_3 \in \bar{X}_3 : \text{L3min} \leq \text{L3} \leq \text{L3max}, \text{P3} = 1\},$$

$$C_{31}^b := \{\bar{x}_3 \in \bar{X}_3 : \text{L3} = \text{L3min}, \text{P3} = 0\},$$

$$C_{32}^a := \{\bar{x}_3 \in \bar{X}_3 :$$

$$\text{L3} \in [\text{L3min}, \text{L3max}], \text{MV3} \in \{1, 2, 3\}, \text{P3} = 1\},$$

$$C_{32}^b := \{\bar{x}_3 \in \bar{X}_3 : \text{L3} \leq \text{L3min}, \text{MV3} \in \{1, 2, 3\}, \text{P3} = 0\},$$

$$C_{33}^a := \{\bar{x}_3 \in \bar{X}_3 : \text{L3} \geq \text{L3min}, \text{MV3} \in \{0, 2, 3\}, \text{P3} = 1\},$$

$$C_{33}^b := \{\bar{x}_3 \in \bar{X}_3 : \text{L3} = \text{L3min}, \text{MV3} \in \{0, 2, 3\}, \text{P3} = 0\},$$

$$C_{34}^a := \{\bar{x}_3 \in \bar{X}_3 : \text{L3} \leq \text{L3min}, \text{MV3} \in \{2, 3\}, \text{P3} = 0\},$$

$$C_{34}^b := \{\bar{x}_3 \in \bar{X}_3 : \text{L3} \geq \text{L3min}, \text{MV3} \in \{2, 3\}, \text{P3} = 1\}.$$

Next, we note that

$$\begin{aligned} & (U(\partial K_e) \setminus K_e) \cap \bar{C}_3 = \\ & (B_1 \cap C_{32}^b) \cup (B_2 \cap C_{33}^a) \cup (B_1 \cap C_{34}^a) \cup (B_2 \cap C_{34}^b), \end{aligned}$$

with

$$\begin{aligned} B_1 \cap C_{32}^b &= \{\bar{x}_3 \in \bar{X}_3 : \\ & \text{L3} \in (\text{L3min} - \epsilon, \text{L3min}), \text{MV3} \in \{1, 2, 3\}, \text{P3} = 0\}, \end{aligned}$$

$$\begin{aligned} B_2 \cap C_{33}^a &= \{\bar{x}_3 \in \bar{X}_3 : \\ & \text{L3} \in (\text{L3max} + f(\tau_3, \text{MV3}), \text{L3max} + f(\tau_3, \text{MV3}) + \epsilon), \\ & \text{MV3} \in \{0, 2, 3\}, \text{P3} = 1\}, \end{aligned}$$

$$\begin{aligned} B_1 \cap C_{34}^a &= \{\bar{x}_3 \in \bar{X}_3 : \\ & \text{L3} \in (\text{L3min} - \epsilon, \text{L3min}), \text{MV3} \in \{2, 3\}, \text{P3} = 0\}, \end{aligned}$$

$$\begin{aligned} B_2 \cap C_{34}^b &= \{\bar{x}_3 \in \bar{X}_3 : \\ & \text{L3} \in (\text{L3max} + f(\tau_3, \text{MV3}), \text{L3max} + f(\tau_3, \text{MV3}) + \epsilon), \\ & \text{MV3} \in \{2, 3\}, \text{P3} = 1\}. \end{aligned}$$

Next, we evaluate the product  $\langle \nabla B(\bar{x}_3), \bar{F}_3(\bar{x}_3, u_3) \rangle$  at each  $(\bar{x}_3, u_3) \in ((U(\partial K_e) \setminus K_e) \cap \bar{C}_3) \times \{0, 1\}$ . Note that

$$\begin{aligned} \nabla B(\bar{x}_3) &= \\ & \begin{bmatrix} 2\text{L3} - (\text{L3max} + \text{L3min}) - f(\tau_3, \text{MV3}) \\ -\text{P3} * (F_{\text{L3}}(3, 0, 1) + \sigma)(\text{L3} - \text{L3min}) \\ \star \\ \star \end{bmatrix}. \end{aligned}$$

Hence,

$$\begin{aligned} \langle \nabla B(\bar{x}_3), \bar{F}_3(\bar{x}_3, u_3) \rangle &= \\ & F_{\text{L3}}(\text{MV3}, \text{P3}, \text{P2}^a)(\text{L3} - \text{L3max} - f(\tau_3, \text{MV3})) + \\ & (\text{L3} - \text{L3min})(F_{\text{L3}}(\text{MV3}, \text{P3}, \text{P2}^a) - \text{P3} * [F_{\text{L3}}(3, 0, 1) + \sigma]). \end{aligned}$$

Next, we distinguish the following two situations:

- 1) When  $\bar{x}_3 \in (B_1 \cap C_{32}^b) \cup (B_1 \cap C_{34}^a)$ , we conclude that  $\text{P3} = 0$ ,  $|\text{L3} - \text{L3min}| \leq \epsilon$ ,  $\text{L3max} - \text{L3} > \text{L3max} - \text{L3min}$ , and

$$F_{\text{L3}}(\text{MV3}, 0, \text{P2}^a) \geq 0 \quad \forall (\text{MV3}, \text{P2}^a) \in \{0, 1, 2\} \times \{0, 1\}.$$

Hence,

$$\begin{aligned} \langle \nabla B(\bar{x}_3), \bar{F}_3(\bar{x}_3, u_3) \rangle &\leq F_{\text{L3}}(\text{MV3}, 0, \text{P2}^a)(\text{L3} - \text{L3max}) \\ &+ \epsilon |F_{\text{L3}}(\text{MV3}, 0, \text{P2}^a)| \\ &\leq -|F_{\text{L3}}(\bar{x}_3, 0, \text{P2}^a)| * (\text{L3max} - \text{L3min} - \epsilon). \end{aligned}$$

Hence, for  $\epsilon$  sufficiently small, we conclude that, for each  $(\bar{x}_3, u_3) \in [(B_1 \cap C_{32}^b) \cup (B_1 \cap C_{34}^a)] \times \{0, 1\}$ ,

$$\langle \nabla B(\bar{x}_3), \bar{F}_3(\bar{x}_3, u_3) \rangle \leq 0.$$

- 2) When  $\bar{x}_3 \in (B_2 \cap C_{33}^a) \cup (B_2 \cap C_{34}^b)$ , we conclude that  $\text{P3} = 1$ ,  $|\text{L3} - \text{L3max} - f(\tau_3, \text{MV3})| \leq \epsilon$ ,  $\text{L3} - \text{L3min} \geq \text{L3max} - \text{L3min}$ , and

$$F_{\text{L3}}(\text{MV3}, 1, \text{P2}^a) - F_{\text{L3}}(3, 0, 1) - \sigma \leq -\sigma$$

for all  $(\text{MV3}, \text{P2}^a) \in \{0, 2, 3\} \times \{0, 1\}$ . Hence,

$$\begin{aligned} \langle \nabla B(\bar{x}_3), \bar{F}_3(\bar{x}_3, u_3) \rangle &\leq |F_{\text{L3}}(\text{MV3}, 1, \text{P2}^a)|\epsilon + \\ &(\text{L3} - \text{L3min})[F_{\text{L3}}(\text{MV3}, 1, \text{P2}^a) - (F_{\text{L3}}(3, 0, 1) + \sigma)] \\ &\leq |F_{\text{L3}}(\text{MV3}, 1, \text{P2}^a)| * \epsilon - \sigma * (\text{L3max} - \text{L3min}). \end{aligned}$$

Hence, for  $\epsilon$  sufficiently small, we conclude that, for each  $(\bar{x}_3, u_3) \in [(B_2 \cap C_{33}) \cup (B_2 \cap C_{34})] \times \{0, 1\}$ ,

$$\langle \nabla B(\bar{x}_3), \bar{F}_3(\bar{x}_3, u_3) \rangle \leq 0.$$

■

So far, we showed that our control logic modification makes the system safe against attacks in P2, but in experimental results our change is not enough when the adversary attacks P1. As a result, we need to modify the PLC controlling stage 2 as well; i.e., C2.

## F.2. Changing the Control Logic of P3 and P2 to Make the System More Secure

**Claim 6.** When modifying the logic in (16) governing P2 as follows:

$$P2 := \begin{cases} 1 & \text{if } (L3 \leq L3_{\min}, P2 = 0, L2 \geq L2_{\min}) \vee \\ & (MV3 \in \{1, 2\}, P2 = 0, L2 \geq L2_{\min}), \\ 0 & \text{if } (L3 \geq L3_{\max}, P2 = 1) \vee \\ & (MV3 \in \{0, 3\}, P2 = 1) \vee \\ & (L2 \leq L2_{\min}, P2 = 1), \end{cases} \quad (87)$$

the plant becomes safe under any attack affecting P1<sup>a</sup>.

•

Since only P1<sup>a</sup> is attacked, the model of Stage 3 is as in (25) and its safety is already analyzed in Theorem 6. Hence, to prove Claim 6, it is enough to prove that Stage 2 is safe uniformly in  $(P1^a, x_3) \in \{0, 1\} \times X_{s3}$  when only P1<sup>a</sup> is attacked and when (87) governs P2. To simplify the proof, we model Stage 2 when only P1<sup>a</sup> is attacked and when the logic governing P2 is modified;

$$\mathcal{H}_2 : \begin{cases} \dot{x}_2 = F_2(x_2, u_2) & (x_2, u_2) \in \tilde{C}_2 \times \{0, 1\} \\ x_2^+ = \tilde{G}_2(x_2, L3, MV3) & (x_2, u_2) \in \tilde{D}_2 \times \{0, 1\}, \end{cases} \quad (88)$$

where

$$u_2 := (L3, MV3, P1^a) \in \mathcal{U}_2,$$

$$\mathcal{U}_2 := [L3_{\min}, L3_{\max} + \delta] \times \{0, 1, 2, 3\} \times \{0, 1\},$$

and

$$F_2(x_2, u_2) := \begin{bmatrix} F_{L2}(MV2, P2, P1^a) \\ 1 \\ 0 \\ 0 \end{bmatrix},$$

**Theorem 13.** Consider the hybrid system  $\mathcal{H}_2$  in (88). Consider the initial set  $X_{o2}$  in (30) and the unsafe set  $X_{u2} \subset X_2 \setminus X_{s2}$  with  $X_{s2}$  introduced in (34). Assume that there exists  $\sigma_h > 0$  such that

$$4T_2(F_{L2}(2, 0, 1) + \sigma_h) \leq \delta. \quad (89)$$

Then, the hybrid system  $\mathcal{H}_2$  in (88) is safe with respect to  $(X_{o2}, X_{u2})$  uniformly in  $u_2 \in \mathcal{U}_2$ , and admits a barrier function certificate given by

$$B(x_2) := (L2 - L2_{\min})(L2 - L2_{\max} - \chi(L2) * h(\tau_2, MV2)), \quad (90)$$

where  $\chi : \mathbb{R} \rightarrow [0, 1]$  is a smooth function such that

$$\begin{cases} \chi(L2) = 1 & \text{if } L2 \geq L2_{\max} \\ \chi(L2) = 0 & \text{if } L2 \leq L2_{\min} \\ \chi(L2) \in [0, 1] & \text{otherwise,} \end{cases}$$

and  $h(\tau_2, MV2) := (F_{L2}(2, 0, 1) + \sigma_h) * [\tau_2 + T_2 * w_h(MV2)]$ ,  $w_h(2) := 0$ ,  $w_h(1) := 1$ ,  $w_h(3) := 2$ ,  $w_h(0) := 3$ . □

*Proof.* Consider the hybrid system  $\mathcal{H}_2$  in (79) with  $u_2 \in \mathcal{U}_2$ . Consider the sets  $(X_{o2}, X_{u2})$  and the barrier function candidate in (90). Note that

$$h(\tau_2, MV2) \in [0, ([F_{L2}(2, 0, 1) + \sigma_h] * T_2 * 4)]$$

for all  $(\tau_2, MV2) \in [0, T_2] \times \{0, 1, 2, 3\}$ . Hence, for  $\sigma_h$  satisfying (83), we conclude that

$$h(\tau_2, MV2) \in [0, \delta] \quad \forall (\tau_2, MV2) \in [0, T_2] \times \{0, 1, 2, 3\}.$$

Thus, (4) is satisfied.

Next, we introduce the notation

$$\bar{u}_2 := (L3, MV3) \in \bar{\mathcal{U}}_2,$$

$$\bar{\mathcal{U}}_2 := [L3_{\min}, L3_{\max} + \delta] \times \{0, 1, 2, 3\}.$$

We also introduce

$$K_e := \{x_2 \in X_2 : L2 \in [L2_{\min}, L2_{\max} + h(\tau_2, MV2)]\} \times \bar{\mathcal{U}}_2.$$

Note that

$$K_w = K_e \times \{0, 1\}.$$

To apply Theorem 1, we start verifying (11) and (12).

Note that the set  $K_e \cap \tilde{D}_2$  satisfies

$$K_e \cap \tilde{D}_2 = A_1 \cup A_2 \cup A_3 \cup A_4 \cup A_5 \cup A_6,$$

where

$$A_1 := \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 :$$

$$L2 \in [L2_{\max}, L2_{\max} + h(\tau_2, 1)], MV2 = 1\},$$

$$A_2 := \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L2 = L2_{\min}, MV2 = 0\},$$

$$A_3 := \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 :$$

$$L2 \in [L2_{\min}, L2_{\max} + h(T_2, MV2)],$$

$$\tau_2 = T_2, MV2 \in \{2, 3\}\},$$

$$A_4 := \{(x_2, \bar{u}_2) \in (X_2 \times \bar{\mathcal{U}}_2) \setminus (A_1 \cup A_2 \cup A_3) :$$

$$L2 \in [L2_{\min}, L2_{\max} + h(\tau_2, MV2)],$$

$$P2 = 1, L3 \geq L3_{\max}\} \cup$$

$$\{(x_2, \bar{u}_2) \in (X_2 \times \bar{\mathcal{U}}_2) \setminus (A_1 \cup A_2 \cup A_3) :$$

$$L2 \in [L2_{\min}, L2_{\max} + h(\tau_2, MV2)],$$

$$P2 = 1, MV3 \in \{0, 3\}\},$$

$$A_5 := \{(x_2, \bar{u}_2) \in (X_2 \times \bar{\mathcal{U}}_2) \setminus (A_1 \cup A_2 \cup A_3) :$$

$$L2 \in [L2_{\min}, L2_{\max} + h(\tau_2, MV2)],$$

$$P2 = 0, L3 = L3_{\min}\} \cup$$

$$\{(x_2, \bar{u}_2) \in (X_2 \times \bar{\mathcal{U}}_2) \setminus (A_1 \cup A_2 \cup A_3) :$$

$$L2 \in [L2_{\min}, L2_{\max} + h(\tau_2, MV2)],$$

$$P2 = 0, MV3 \in \{1, 2\}\},$$

$$A_6 := \{(x_2, \bar{u}_2) \in (X_2 \times \bar{\mathcal{U}}_2) \setminus (A_1 \cup A_2 \cup A_3) :$$

$$L2 = L2_{\min}, P2 = 1\}.$$

Note that, for each  $(x_2, \bar{u}_2) \in A_1$ , we have

$$G_2(x_2, \bar{u}_2) \in [L2 \ 0 \ 3 \ \{0, 1\}]^\top.$$

Hence,

$$\begin{aligned} B(G_2(x_2, \bar{u}_2)) &= (L_2 - L_{2\min}) * \\ &\quad (L_2 - L_{2\max} - \chi(L_2) * h(0, 3)) \\ &\leq 0. \end{aligned}$$

The latter inequality is true since

$$L_2 \geq L_{2\max} \quad \forall (x_2, \bar{u}_2) \in A_1,$$

and

$$h(\tau_2, 1) \leq h(0, 3) \quad \forall \tau_2 \in [0, T_2].$$

Similarly, for each  $(x_2, \bar{u}_2) \in A_2$ , we have

$$G_2(x_2, \bar{u}_2) = [L_2 \ 0 \ 2 \ \{0, 1\}]^\top.$$

Hence,

$$\begin{aligned} B(G_2(x_2, \bar{u}_2)) &= (L_2 - L_{2\min}) * \\ &\quad (L_2 - L_{2\max} - \chi(L_2) * h(0, 2)) \\ &\leq 0. \end{aligned}$$

The latter inequality is true since  $h(0, 2) \geq 0$  and

$$L_2 \leq L_{2\min} \quad \forall (x_2, \bar{u}_2) \in A_2.$$

Now, for each  $(x_2, \bar{u}_2) \in A_3$ ,

$$G_2(x_2, \bar{u}_2) \in [L_2 \ 0 \ \alpha(MV_2) \ \{0, 1\}]^\top,$$

where  $\alpha(3) := 0$  and  $\alpha(2) := 1$ . Hence,

$$\begin{aligned} B(G_2(x_2, \bar{u}_2)) &= (L_2 - L_{2\min}) * \\ &\quad (L_2 - L_{2\max} - \chi(L_2) * h(0, \alpha(MV_2))) \leq 0. \end{aligned}$$

The latter inequality is true since

$$h(7, 3) \leq h(0, 0), \quad h(7, 2) \leq h(0, 1).$$

Next, for each  $(x_2, \bar{u}_2) \in A_4 \cup A_5$ , we have

$$G_2(x_2, \bar{u}_2) \in [L_2 \ \tau_2 \ MV_2 \ \{0, 1\}]^\top.$$

Hence,

$$\begin{aligned} B(G_2(x_2, \bar{u}_2)) &= (L_2 - L_{2\min}) * \\ &\quad (L_2 - L_{2\max} - \chi(L_2) * h(\tau_2, MV_2)) \leq 0. \end{aligned}$$

The latter inequality is true since  $(A_4 \cup A_5) \subset K_e$ . Hence, we conclude that (11) is satisfied. Moreover, to show (12), we note that  $\tilde{C}_2 \cup \tilde{D}_2 = X_2 \times \bar{U}_2$  and  $G_2(x_2, \bar{u}_2) \in X_2$  for all  $(x_2, \bar{u}_2) \in \tilde{D}_2$ .

Next, to verify (10), we start noting that the set  $U(\partial K_e) \setminus K_e$  satisfies

$$U(\partial K_e) \setminus K_e = B_1 \cup B_2,$$

where

$$B_1 := \{(x_2, \bar{u}_2) \in X_2 \times \bar{U}_2 : L_2 \in (L_{2\min} - \epsilon, L_{2\min})\}$$

and

$$\begin{aligned} B_2 := \{(x_2, \bar{u}_2) \in X_2 \times \bar{U}_2 : \\ L_2 \in (L_{2\max} + h(\tau_2, MV_2), L_{2\max} + h(\tau_2, MV_2) + \epsilon)\}. \end{aligned}$$

Furthermore, to compute the set  $U(\partial K_e) \setminus K_e \cap \tilde{C}_2$ , we compute the set  $\tilde{C}_2$  using the following Lemma.

**Lemma 1.** The flow set  $\tilde{C}_2$  is given by

$$\tilde{C}_2 = [C_2 \cap K_1] \cup [D_{23} \cap K_2],$$

where

$$\begin{aligned} K_1 &:= \text{cl}((X_2 \times \bar{U}_2) \setminus \tilde{D}_{26}), \quad K_2 := \text{cl}((X_2 \times \bar{U}_2) \setminus \tilde{I}), \\ \tilde{I} &:= \tilde{D}_{26} \cup \{(x_2, \bar{u}_2) \in X_2 \times \bar{U}_2 : P_2 = 0, L_2 \geq L_{2\min}\}, \\ \tilde{D}_{26} &:= \{(x_2, \bar{u}_2) \in X_2 \times \bar{U}_2 : L_2 \leq L_{2\min}, P_2 = 1\}, \\ D_{23} &:= \{(x_2, \bar{u}_2) \in X_2 \times \bar{U}_2 : L_3 \leq L_{3\min}, P_2 = 0\} \cup \\ &\quad \{(x_2, \bar{u}_2) \in X_2 \times \bar{U}_2 : MV_3 \in \{1, 2\}, P_2 = 0\}. \end{aligned}$$

□

*Proof.* Note that the sets  $D_2$  and  $\tilde{D}_2$  can be expressed as follows:

$$\begin{aligned} D_2 &:= D_{21} \cup D_{22} \cup D_{23} \cup D_{24} \cup D_{25}, \\ \tilde{D}_2 &:= D_{21} \cup D_{22} \cup \tilde{D}_{23} \cup D_{24} \cup D_{25} \cup \tilde{D}_{26}, \end{aligned}$$

with

$$\begin{aligned} D_{21} &:= \{(x_2, \bar{u}_2) \in X_2 \times \bar{U}_2 : L_2 \geq L_{2\max}, MV_2 = 1\} \\ D_{22} &:= \{(x_2, \bar{u}_2) \in X_2 \times \bar{U}_2 : \tau_2 \geq T_2, MV_2 \in \{2, 3\}\} \\ D_{24} &:= \{(x_2, \bar{u}_2) \in X_2 \times \bar{U}_2 : L_2 \leq L_{2\min}, MV_2 = 0\} \\ D_{25} &:= \{(x_2, \bar{u}_2) \in X_2 \times \bar{U}_2 : L_3 \geq L_{3\max}, P_2 = 1\} \cup \\ &\quad \{(x_2, \bar{u}_2) \in X_2 \times \bar{U}_2 : MV_3 \in \{0, 3\}, P_2 = 1\}. \\ \tilde{D}_{23} &:= \{(x_2, \bar{u}_2) \in X_2 \times \bar{U}_2 : L_3 \leq L_{3\min}, P_2 = 0, \\ &\quad L_2 \geq L_{2\min}\} \cup \{(x_2, \bar{u}_2) \in X_2 \times \bar{U}_2 : MV_3 \in \{1, 2\}, \\ &\quad P_2 = 0, L_2 \geq L_{2\min}\}. \end{aligned}$$

Furthermore, note that  $\tilde{D}_{23} = D_{23} \cap I$ , where

$$I := \{(x_2, \bar{u}_2) \in X_2 \times \bar{U}_2 : P_2 = 0, L_2 \geq L_{2\min}\}.$$

Hence, we have that

$$\begin{aligned} \tilde{D}_2 &= D_{21} \cup D_{22} \cup (D_{23} \cap I) \cup D_{24} \cup D_{25} \cup \tilde{D}_{26} \\ &= [D_{21} \cup D_{22} \cup D_{24} \cup D_{25}] \cup (D_{23} \cap I) \cup \tilde{D}_{26} \\ &= \{(D_{21} \cup D_{22} \cup D_{24} \cup D_{25} \cup D_{23}) \cap \\ &\quad ([D_{21} \cup D_{22} \cup D_{24} \cup D_{25}] \cup I)\} \cup \tilde{D}_{26} \\ &= \{D_2 \cap ([D_2 \setminus D_{23}] \cup I)\} \cup \tilde{D}_{26} \\ &= (D_2 \cup \tilde{D}_{26}) \cap (([D_2 \setminus D_{23}] \cup I) \cup \tilde{D}_{26}) \\ &= (D_2 \cup \tilde{D}_{26}) \cap ([D_2 \setminus D_{23}] \cup I \cup \tilde{D}_{26}) \\ &= (D_2 \cup \tilde{D}_{26}) \cap ([D_2 \setminus D_{23}] \cup \tilde{I}). \end{aligned}$$

Next, we note that

$$\tilde{I} := \tilde{D}_{26} \cup \{(x_2, \bar{u}_2) \in X_2 \times \bar{U}_2 : P_2 = 0, L_2 \geq L_{2\min}\}.$$

Hence,

$$\begin{aligned} \tilde{C}_2 &= \text{cl}[(X_2 \times \bar{U}_2) \setminus \tilde{D}_2] \\ &= \text{cl}[(X_2 \times \bar{U}_2) \setminus ((D_2 \cup \tilde{D}_{26}) \cap ([D_2 \setminus D_{23}] \cup \tilde{I}))] \\ &= (C_2 \cap \text{cl}[(X_2 \times \bar{U}_2) \setminus \tilde{D}_{26}]) \cup \\ &\quad \{(D_{23} \cup C_2) \cap \text{cl}[(X_2 \times \bar{U}_2) \setminus \tilde{I}]\} \\ &= (C_2 \cap \text{cl}[(X_2 \times \bar{U}_2) \setminus \tilde{D}_{26}]) \cup \{D_{23} \cap \text{cl}[(X_2 \times \bar{U}_2) \setminus \tilde{I}]\} \\ &= (C_2 \cap K_1) \cup \{D_{23} \cap K_2\}. \end{aligned}$$

■

Note that

$$K_1 = K_2 \cup \{(x_2, \bar{u}_2) \in X_2 \times \bar{U}_2 : L_2 \geq L_{2\min}, P_2 = 0\},$$

and

$$K_2 = \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_2 \geq L_{2\min}, P_2 = 1\} \cup \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_2 \leq L_{2\min}, P_2 = 0\}.$$

Hence,

$$D_{23} \cap K_2 := \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_3 \leq L_{3\min}, P_2 = 0, L_2 \leq L_{2\min}\} \cup \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : MV_3 \in \{1, 2\}, P_2 = 0, L_2 \leq L_{2\min}\}.$$

Next, we compute the set  $C_2 \cap K_1$ , which can be expressed as

$$C_2 \cap K_1 = \bigcup_{i=1}^8 \tilde{C}_{2i},$$

where

$$\tilde{C}_{2i} := \tilde{C}_{2i}^a \cup \tilde{C}_{2i}^b \quad \forall i \in \{2, 4, 6, 8\},$$

and

$$\tilde{C}_{21} := \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_2 \leq L_{2\max}, MV_2 = \{1, 2, 3\}, P_2 = \{0\}, MV_3 = \{0, 3\}\},$$

$$\tilde{C}_{22}^a := \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_2 \leq L_{2\max}, MV_2 = \{1, 2, 3\}, L_3 \leq L_{3\max}, P_2 = \{0\}, MV_3 = \{0, 3\}\},$$

$$\tilde{C}_{22}^b := \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_2 \in [L_{2\min}, L_{2\max}], MV_2 = \{1, 2, 3\}, L_3 \leq L_{3\max}, P_2 = \{1\}, MV_3 = \{1, 2\}\},$$

$$\tilde{C}_{23} := \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_{2\min} \leq L_2 \leq L_{2\max}, P_2 = \{0\}, MV_3 = \{0, 3\}\},$$

$$\tilde{C}_{24}^a := \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_{2\min} \leq L_2 \leq L_{2\max}, L_3 \leq L_{3\max}, P_2 = \{0\}, MV_3 = \{0, 3\}\},$$

$$\tilde{C}_{24}^b := \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_{2\min} \leq L_2 \leq L_{2\max}, L_3 \leq L_{3\max}, P_2 = \{1\}, MV_3 = \{1, 2\}\},$$

$$\tilde{C}_{25} := \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : MV_2 \in \{2, 3\}, P_2 = \{0\}, MV_3 = \{0, 3\}\},$$

$$\tilde{C}_{26}^a := \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : MV_2 \in \{2, 3\}, L_3 \leq L_{3\max}, P_2 = \{0\}, MV_3 = \{0, 3\}\},$$

$$\tilde{C}_{26}^b := \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_2 \geq L_{2\min}, MV_2 \in \{2, 3\}, L_3 \leq L_{3\max}, P_2 = \{1\}, MV_3 = \{1, 2\}\},$$

$$\tilde{C}_{27} := \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_2 \geq L_{2\min}, MV_2 = \{0, 2, 3\}, P_2 = \{0\}, MV_3 = \{0, 3\}\},$$

$$\tilde{C}_{28}^a := \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_2 \geq L_{2\min}, MV_2 = \{0, 2, 3\}, L_3 \in [L_{3\min}, L_{3\max}], P_2 = \{0\}, MV_3 = \{0, 3\}\}.$$

$$\tilde{C}_{28}^b := \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_2 \geq L_{2\min}, MV_2 = \{0, 2, 3\}, L_3 \in [L_{3\min}, L_{3\max}], P_2 = \{1\}, MV_3 = \{1, 2\}\}.$$

Hence, we conclude that

$$(U(\partial K_e) \setminus K_e) \cap \tilde{C}_2 = (\tilde{C}_{21} \cap B_1) \cup (\tilde{C}_{22}^a \cap B_1) \cup (\tilde{C}_{25} \cap B_1) \cup ((D_{23} \cap K_2) \cap B_1) \cup (\tilde{C}_{26} \cap B_1) \cup (\tilde{C}_{25} \cap B_2) \cup (\tilde{C}_{26} \cap B_2) \cup (\tilde{C}_{27} \cap B_2) \cup (\tilde{C}_{28} \cap B_2),$$

with

$$\tilde{C}_{21} \cap B_1 = \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_2 \in (L_{2\min} - \epsilon, L_{2\min}), MV_2 = \{1, 2, 3\}, P_2 = \{0\}, MV_3 = \{0, 3\}\},$$

$$\tilde{C}_{22} \cap B_1 = \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_2 \in (L_{2\min} - \epsilon, L_{2\min}), MV_2 = \{1, 2, 3\}, L_3 \leq L_{3\max}, P_2 = \{0\}, MV_3 = \{0, 3\}\},$$

$$\tilde{C}_{25} \cap B_1 = \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_2 \in (L_{2\min} - \epsilon, L_{2\min}), MV_2 \in \{2, 3\}, P_2 = \{0\}, MV_3 = \{0, 3\}\},$$

$$\tilde{C}_{26} \cap B_1 = \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_2 \in (L_{2\min} - \epsilon, L_{2\min}), MV_2 \in \{2, 3\}, L_3 \leq L_{3\max}, P_2 = \{0\}, MV_3 = \{0, 3\}\},$$

$$((D_{23} \cap K_2) \cap B_1) = \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_3 \leq L_{3\min}, P_2 = 0, L_2 \in (L_{2\min} - \epsilon, L_{2\min})\} \cup \{(x_2, u_2) \in X_2 \times \mathcal{U}_2 : MV_3 \in \{1, 2\}, P_2 = 0, L_2 \in (L_{2\min} - \epsilon, L_{2\min})\},$$

$$\tilde{C}_{25} \cap B_2 = \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_2 \in (L_{2\max} + h(\tau_2, MV_2), L_{2\max} + h(\tau_2, MV_2) + \epsilon), MV_2 \in \{2, 3\}, P_2 = \{0\}, MV_3 = \{0, 3\}\},$$

$$\tilde{C}_{26} \cap B_2 = \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_2 \in (L_{2\max} + h(\tau_2, MV_2), L_{2\max} + h(\tau_2, MV_2) + \epsilon), MV_2 \in \{2, 3\}, L_3 \leq L_{3\max}, P_2 = \{0\}, MV_3 = \{0, 3\}\} \cup \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_2 \in (L_{2\max} + h(\tau_2, MV_2), L_{2\max} + h(\tau_2, MV_2) + \epsilon), MV_2 \in \{2, 3\}, L_3 \leq L_{3\max}, P_2 = \{1\}, MV_3 = \{1, 2\}\},$$

$$\tilde{C}_{27} \cap B_2 = \{(x_2, \bar{u}_2) \in X_2 \times \bar{\mathcal{U}}_2 : L_2 \in (L_{2\max} + h(\tau_2, MV_2), L_{2\max} + h(\tau_2, MV_2) + \epsilon), MV_2 = \{0, 2, 3\}, P_2 = \{0\}, MV_3 = \{0, 3\}\},$$

$$\tilde{C}_{28} \cap B_2 = \{(x_2, \bar{u}_2) \in X_2 \times \bar{U}_2 :$$

$$\begin{aligned} &L_2 \in (L_{2\max} + h(\tau_2, MV_2), L_{2\max} + h(\tau_2, MV_2) + \epsilon), \\ &MV_2 = \{0, 2, 3\}, L_3 \in [L_{3\min}, L_{3\max}], P_2 = \{1\}, \\ &MV_3 = \{1, 2\} \cup \{(x_2, \bar{u}_2) \in X_2 \times \bar{U}_2 : \\ &L_2 \in (L_{2\max} + h(\tau_2, MV_2), L_{2\max} + h(\tau_2, MV_2) + \epsilon), \\ &MV_2 = \{0, 2, 3\}, L_3 \in [L_{3\min}, L_{3\max}], P_2 = \{0\}, \\ &MV_3 = \{0, 3\}\}. \end{aligned}$$

Next, we evaluate the product  $\langle \nabla B(x_2), F_2(x_2, u_2) \rangle$  at each  $(x_2, u_2) \in ((U(\partial K_e) \setminus K_e) \cap \tilde{C}_2) \times \{0, 1\}$ . Note that

$$\nabla B(x_2) = [\nabla_1 B(x_2) \ \nabla_2 B(x_2) \ \star \ \star]^\top,$$

where

$$\begin{aligned} \nabla_1 B(x_2) := & 2L_2 - (L_{2\max} + L_{2\min}) - \\ & \chi(L_2) * h(\tau_2, MV_2) - \\ & (L_2 - L_{2\min}) * h(\tau_2, MV_2) * \frac{\partial \chi}{\partial L_2}(L_2), \end{aligned}$$

$$\nabla_2 B(x_2) := -[L_2 - L_{2\min}] * \chi(L_2) * [F_{L_2}(2, 0, 1) + \sigma_h].$$

Hence,

$$\begin{aligned} \langle \nabla B(x_2), F(x_2, u_2) \rangle = & F_{L_2}(MV_2, P_2, P_1^a) * \\ & [2L_2 - (L_{2\min} + L_{2\max}) - \chi(L_2) * h(\tau_2, MV_2)] - \\ & F_{L_2}(MV_2, P_2, P_1^a) * \\ & (L_2 - L_{2\min}) * h(\tau_2, MV_2) * \frac{\partial \chi}{\partial L_2}(L_2) - \\ & [L_2 - L_{2\min}] * \chi(L_2) * [F_{L_2}(2, 0, 1) + \sigma_h]. \end{aligned}$$

Next, we distinguish the following two situations:

- 1) When  $(x_2, \bar{u}_2) \in (\tilde{C}_{21} \cap B_1) \cup (\tilde{C}_{22}^a \cap B_1) \cup (\tilde{C}_{25} \cap B_1) \cup ((D_{23} \cap K_2) \cap B_1) \cup (\tilde{C}_{26} \cap B_1)$ , we conclude that  $P_2 = 0$ ,  $\chi(L_2) = 0$ , and  $\frac{\partial \chi}{\partial L_2}(L_2) = 0$ . Hence,
$$\langle \nabla B(x_2), F(x_2, u_2) \rangle = -F_{L_2}(MV_2, 0, P_1^a) * [L_{2\min} + L_{2\max} - 2L_2] \leq 0.$$

The latter inequality is true since according to (77), we have

$$\begin{aligned} F_{L_2}(MV_2, 0, P_1^a) &\geq 0 \\ \forall (MV_2, P_1^a) \in \{0, 1, 2\} \times \{0, 1\}. \end{aligned}$$

- 2) When  $(x_2, u_2) \in (\tilde{C}_{25} \cap B_2) \cup (\tilde{C}_{26} \cap B_2) \cup (\tilde{C}_{27} \cap B_2) \cup (\tilde{C}_{28} \cap B_2)$ , we conclude that  $\chi(L_2) = 1$ , and  $\frac{\partial \chi}{\partial L_2}(L_2) = 0$ . Hence,

$$\begin{aligned} F_{L_2}(MV_2, P_2, P_1^a) &\leq F_{L_2}(2, 0, 1), \\ 0 \leq L_2 - L_{2\max} - h(\tau_2, MV_3) &\leq \epsilon. \end{aligned}$$

Hence,

$$\begin{aligned} \langle \nabla B(x_2), F_2(x_2, u_2) \rangle &\leq F_{L_2}(MV_2, P_2, P_1^a) * \\ &[L_2 - L_{2\min} + \epsilon] - \\ &[L_2 - L_{2\min}] * [F_{L_2}(2, 0, 1) + \sigma_h] \end{aligned}$$

and, for  $\epsilon$  sufficiently small, we obtain

$$\langle \nabla B(x_2), F_2(x_2, u_2) \rangle \leq -\sigma_h[L_{2\max} - L_{2\min}]/2.$$

Hence,

$$\langle \nabla B(x_2), F_2(x_2, u_2) \rangle \leq 0.$$

■

Finally, When the Attacker Enforces  $MV_1^a = 1$ , in this case, the water level in the first tank will overflow. Indeed, the quantity of water added to the first tank in this case is bigger than the quantity removed, even when  $P_1^a = 1$  and  $MV_2^a = 1$ .

When the Attacker Enforces  $MV_1^a = 0$ , in this case, starvation of the first stage occurs due to the constant demand of water when  $P_3 = 1$ . However, this unsafe behavior can be avoided if

- a. we use the logic in (84) governing  $P_3$  with  $L_{3o} = L_{3\min}$ .
- b. we modify the logic in (16) governing  $P_2$  using (87).
- c. we modify the logic governing  $P_1$  using

$$P_1 := \begin{cases} 1 & \text{if } (L_2 \leq L_{2\min}, P_1 = 0, L_1 \geq L_{1\min}) \vee \\ & (MV_2 \in \{1, 2\}, P_1 = 0, L_2 \geq L_{1\min}), \\ 0 & \text{if } (L_2 \geq L_{2\max}, P_1 = 1) \vee \\ & (MV_2 \in \{0, 3\}, P_1 = 1) \vee (L_2 \leq L_{1\min}, P_1 = 1). \end{cases} \quad (91)$$

If the logic governing one of the pumps ( $P_1, P_2, P_3$ ) is not modified according to the aforementioned items, then the corresponding stage will be subject to starvation.

When the Attacker Targets  $MV_2^a$  and  $MV_3^a$  Only and Enforces  $MV_2^a = 1$  or  $MV_3^a = 1$ , in this case, the system remains safe since the pumps prevent overflows when they are closed.

When the Attacker Targets  $MV_2^a$  Only and Enforces  $MV_2^a = 0$ , in this case, starvation in both Stages 2 and 3 occurs. However, if we use the logic in (84) governing  $P_3$  with  $L_{3o} = L_{3\min}$ , we avoid starvation of Stage 3. Moreover, if we additionally modify the logic in (16) governing  $P_2$  using (87), the system becomes safe.

When the Attacker Targets  $MV_3^a$  Only and Enforces  $MV_3^a = 0$ , in this case, starvation in Stage 3 occurs. However, if we use the logic in (84) governing  $P_3$  with  $L_{3o} = L_{3\min}$ , the system becomes safe.

In summary, we have shown that the original SWaT system cannot guarantee safety when the attacker compromises any of the following actuators  $P_1, P_2, MV_1, MV_2, MV_3$ . However, we proposed a set of control logic changes to PLCs 1 and 2 and with these changes we were able to prove that the system is safe if the attacker compromises any of these actuators:  $P_1, P_2, P_3, MV_2, MV_3$ . The only time SWaT cannot guarantee safety is when the attacker compromises  $MV_1$ . The reason for this is that the amount of water coming into the first tank is higher than the amount that can be taken out by  $P_1$ . To guarantee safety against a compromise of  $MV_1$  we would need a physical redesign of the system so that the rate of flow of entering water is the same as the rate of flow that  $P_1$  can take out of the first tank.