# A Virtual Cybersecurity Department for Securing Digital Twins in Water Distribution Systems

1st Mohammadhossein Homaei
*Media Engineering Group*
*University of Extremadura*
Cáceres, Spain
mhomaein@alumnos.unex.es

2nd Agustin Di Bartolo
3rd Óscar Mogollón-Gutiérrez
*Media Engineering Group*
*University of Extremadura*
Cáceres, Spain
{adibartolo, oscarmg}@unex.es

4th Fernando Broncano Morgado,
5th Pablo García Rodríguez
*Media Engineering Group*
*University of Extremadura*
Cáceres, Spain
{fbroncano, pablogr}@unex.es

*Abstract*—**Digital twins (DTs) help improve real-time monitoring and decision-making in water distribution systems. However, their connectivity makes them easy targets for cyberattacks such as scanning, denial-of-service (DoS), and unauthorized access. Small and medium-sized enterprises (SMEs) that manage these systems often do not have enough budget or staff to build strong cybersecurity teams. To solve this problem, we present a Virtual Cybersecurity Department (VCD), an affordable and automated framework designed for SMEs. The VCD uses open-source tools like Zabbix for real-time monitoring, Suricata for network intrusion detection, Fail2Ban to block repeated login attempts, and simple firewall settings. To improve threat detection, we also add a machine-learning-based IDS trained on the OD-IDS2022 dataset using an improved ensemble model. This model detects cyber threats such as brute-force attacks, remote code execution (RCE), and network flooding, with 92% accuracy and fewer false alarms. Our solution gives SMEs a practical and efficient way to secure water systems using low-cost and easy-to-manage tools.**

*Index Terms*—**Digital Twins, Cybersecurity, Intrusion Detection System, Machine Learning, Zabbix, Water Distribution, SMEs**

## I. INTRODUCTION

As water distribution systems become increasingly connected, they face growing cybersecurity risks [1]–[6]. Integrating information technology (IT) and operational technology (OT) has significantly improved efficiency and real-time monitoring, but has also introduced new vulnerabilities. DT technology, which provides virtual replicas of physical systems, further enhances these capabilities by improving operational visibility, predictive maintenance, and decision-making [7], [8]. However, increased connectivity and intelligence in DTs expand their attack surface, making them vulnerable not only to data leaks but also to threats that could impact public health and infrastructure safety. Cyberattacks such as unauthorized access, data manipulation, or DoS could result in severe incidents like water contamination or system disruptions [9]. These attacks can go undetected for long periods, especially in systems without active monitoring. As the number of smart sensors and connected devices grows, the complexity of protecting the infrastructure increases. Thus, robust, automated cybersecurity measures are crucial.

SMEs, which often manage water distribution networks, have serious challenges in cybersecurity because they usually do not have enough money or trained IT staff. Traditional security systems are expensive and need expert teams, so they are not good options for small organizations. To solve this, we propose a VCD, a low-cost and easy-to-use system built with open-source tools. The main tool in the VCD is Zabbix, which gives real-time system monitoring, alerting, and data visualization [**?**], [10]. It helps detect technical problems and possible cyber threats in the digital twin environment. To improve detection, we also added a machine-learning-based IDS trained on the OD-IDS2022 dataset. This system can find different types of cyberattacks, such as scanning, brute-force, RCE, and DoS attacks. By combining simple monitoring with advanced machine learning, our framework gives SMEs an effective and affordable way to protect their water systems.

Unlike many existing frameworks, our VCD uniquely combines traditional open-source tools with a customized, explainable machine learning model, all optimized for low-resource environments typical in SME water utilities.

The motivation for this work comes from real needs in the field. Many small and rural water utilities want to use digital twin systems, but they are not ready to face growing cybersecurity risks. Most existing solutions are made for large companies and need expensive tools or professional IT teams. SMEs cannot afford these systems and are left with weak protection. Also, new cyberattacks are becoming smarter and harder to detect with old methods. Our goal is to offer a practical and affordable solution that helps SMEs protect their water infrastructure using tools they can manage themselves. The proposed Virtual Cybersecurity Department gives them a way to use open-source software, automate responses, and improve detection with machine learning—without needing large investments or complex systems.

The remainder of the paper is organized as follows: Section II reviews existing research in cybersecurity for DT-based water systems. Section III introduces our proposed VCD framework, including system architecture, communication flow, cybersecurity integration, and ML-based IDS. Section IV presents the experimental evaluation, including results from the Zabbix-based monitoring and the IDS model performance. Finally, Section V concludes the paper and outlines directions for future work.

TABLE I
RECENT WORK ON CYBERSECURITY IN DT-ENABLED WATER DISTRIBUTION SYSTEMS (POST-2020)

| Ref. | Focus, Challenges, and Tech. & Eval. | | |
|---|---|---|---|
| | *Focus* | *Challenges* | *Tech. & Eval.* |
| Zhang *et al.* 2021 [12] | Attack detection in DT water systems | Distinguish anomalies from normal ops; Integrate IT/OT data | ML anomaly detection; IoT integration; Simulation testbed |
| Liu *et al.* 2022 [13] | Secure smart water DTs | Ensuring secure communication | Cryptographic protocols; Anomaly-based IDS; Emulated DT with intrusions |
| Qi *et al.* 2022 [14] | Risk assessment in DT networks | Prioritizing vulnerabilities in distributed systems | Sensor fusion; Statistical threat scoring; Risk evaluation scenario analysis |
| Kumar *et al.* 2023 [15] | Mitigate attacks via anomaly+blockchain | Data tampering, traceability | Blockchain for data integrity; ML detection; Experimental deployment |
| Lin *et al.* 2023 [16] | IDS using DT correlation (hydraulic/network) | Detect stealthy attacks in normal ops | Hybrid IDS correlating physical & network metrics; Lab-scale DT with synthetic attacks |

## II. RELATED WORK

### A. Cybersecurity in DTs for Water Systems

In recent years, DT technology has become more common in water distribution systems due to its ability to provide real-time monitoring, predictive analytics, and decision support. However, this increased connectivity has also introduced new cybersecurity challenges. DTs, by design, connect multiple physical and digital components, which increases the attack surface for potential cyber threats.

Zhang *et al.* [12] presented a machine-learning-based intrusion detection framework for DT-enabled water systems, integrating IoT sensors to detect anomalies in physical and cyber operations. Liu *et al.* [13] proposed a secure DT architecture using encryption protocols and anomaly-based IDS to protect communication flows between devices and the cloud. Homaei *et al.* [1], [2] also highlighted the dual role of DTs as both monitoring tools and high-risk targets for attacks, especially in rural water networks.

These studies show that while DTs improve operations, they also require new security solutions that go beyond traditional IT protections.

### B. Challenges in DT-based Water Infrastructure

Cybersecurity in water distribution systems faces multiple technical and operational challenges, particularly when DTs are integrated:

- Anomaly Detection: DT systems rely on normal behavior patterns to function correctly. However, cyberattacks often mimic legitimate fluctuations (e.g., consumption peaks), making detection difficult without advanced ML techniques.
- Scalability and Performance: Real-time monitoring and analysis require high processing power and efficient algorithms, especially as the number of IoT sensors increases.
- Legacy and Modern System Integration: Many utilities still use legacy systems that are not easily compatible with modern IoT devices or secure communication protocols, creating interoperability issues.
- Network Communication Risks: Protocols used in DTs are sometimes unencrypted or misconfigured, exposing them to packet sniffing, spoofing, or DoS attacks.

- Limited Resources in SMEs: Most SMEs lack the IT staff, funding, or training to maintain enterprise-level cybersecurity systems, leaving them especially vulnerable.
- Public Safety and Reliability: Failures in cyber-protected DTs could lead to water shortages, contamination, or service disruptions, affecting entire communities.

These issues make it clear that new frameworks should be lightweight, scalable, and capable of operating in low-resource environments.

### C. Emerging Solutions and Gaps

Recent research has introduced several approaches to improve cybersecurity in DT-enabled water networks. Qi *et al.* [14] introduced a risk assessment method using sensor fusion and statistical analysis to identify vulnerable components. Kumar *et al.* [15] proposed combining blockchain with anomaly detection to increase data traceability and prevent tampering. Lin *et al.* [16] focused on hybrid intrusion detection systems that analyze both physical process data and network logs to detect stealth attacks.

Although these methods show progress, they often rely on complex systems or high-performance resources, which may not be suitable for SMEs.

### D. Positioning of This Work

In contrast to prior works that require extensive infrastructure or expert personnel, our proposed VCD offers a practical alternative for small and medium-sized enterprises. The VCD uses a combination of lightweight, open-source tools—Zabbix, Suricata, and Fail2Ban—alongside a machine-learning-based IDS trained on the OD-IDS2022 dataset.

Unlike many traditional systems that depend solely on signature-based detection or manual log review, our model integrates real-time monitoring with automated responses and a trained ensemble ML model. This hybrid approach improves detection of advanced threats such as brute-force, RCE, and DoS attacks, making it well-suited for decentralized water systems with limited resources. It also reduces the need for continuous human supervision and simplifies system maintenance, allowing operators to focus on operational tasks rather than complex cybersecurity management.
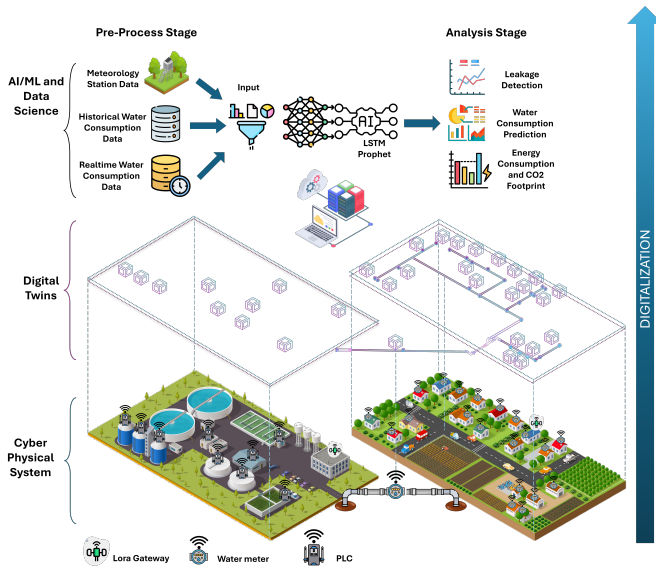
Fig. 1. DT platform in the WDS [17]


Fig. 2. Deployment of Zabbix proxies on Raspberry Pi devices for real-time data collection from IoT meters and SCADA systems.

## III. PROPOSED FRAMEWORK

This section describes the structure and components of the proposed VCD, a cost-effective monitoring framework for DTs in WDS. The system is designed to help SMEs enhance their operational security through automated, open-source tools. The framework includes four main components: the DT system overview, system architecture and communication flow, cyber-security integration using Zabbix and Suricata, and a machine learning-based IDS.

### A. DT System Overview

The VCD is built on a Digital Twin platform that integrates real-time data collection, AI-driven analytics, and secure communication. It consists of three main layers: cyber-physical systems (CPS), data management, and predictive analytics.

The CPS layer includes sensors, PLCs, and IoT water meters deployed in water treatment facilities and distribution pipelines. These devices collect environmental, operational, and consumption data. The data is transmitted securely using technologies such as LoRaWAN, VPN, and SSH. AI/ML models—including LSTM, Prophet, and LightGBM—are used for water usage forecasting, leakage detection, and energy monitoring. Additionally, GIS tools support spatial analysis and map-based monitoring (Figure 1) [17].

This architecture is designed for rural and small-scale water utilities but is scalable for larger infrastructures. It provides enhanced operational control, cost efficiency, and resource optimization.

### B. System Architecture and Communication Flow

The system consists of three key components: edge nodes, secure communication channels, and a central server.

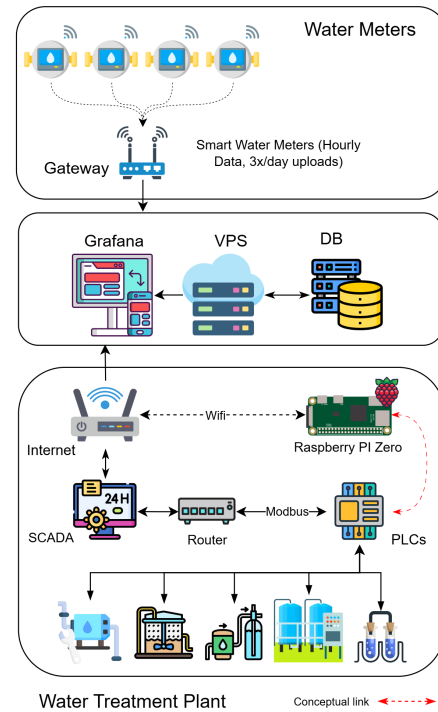- Edge nodes include Raspberry Pi devices equipped with Zabbix proxies, IoT meters, SCADA units, and PLCs. These nodes are strategically placed at water plants and administrative locations to ensure complete visibility (Figure 2).
- Secure communication is established using VPN tunnels, SSH protocols, and LoRaWAN networks. This ensures that data from the field devices reaches the central server with integrity and confidentiality. Zabbix continuously monitors the stability and quality of these connections.
- The central server, hosted on a Virtual Private Server (VPS), aggregates all incoming data. It runs Zabbix for real-time monitoring, Suricata for intrusion detection, and Fail2Ban for automated IP blocking. The server can optionally connect to cloud platforms like AWS or Azure for data storage and computational scalability.

Figure 3 presents the VCD architecture, highlighting the placement of Zabbix and the ML-based IDS modules.

### C. Cybersecurity Integration with Zabbix and Suricata

The core of the cybersecurity layer is Zabbix, which provides data collection, visualization, and alerting functionalities. It monitors metrics such as network traffic, CPU load, memory usage, and failed login attempts. Zabbix is integrated with Suricata, an open-source IDS that inspects network packets and detects threats like port scanning, brute-force logins, and unusual data flows. Suricata's alerts are visualized in the Zabbix dashboard. Fail2Ban complements the system by monitoring authentication logs. It automatically bans IP addresses that exceed a defined number of failed login attempts. This combination of tools ensures multi-layered protection against a
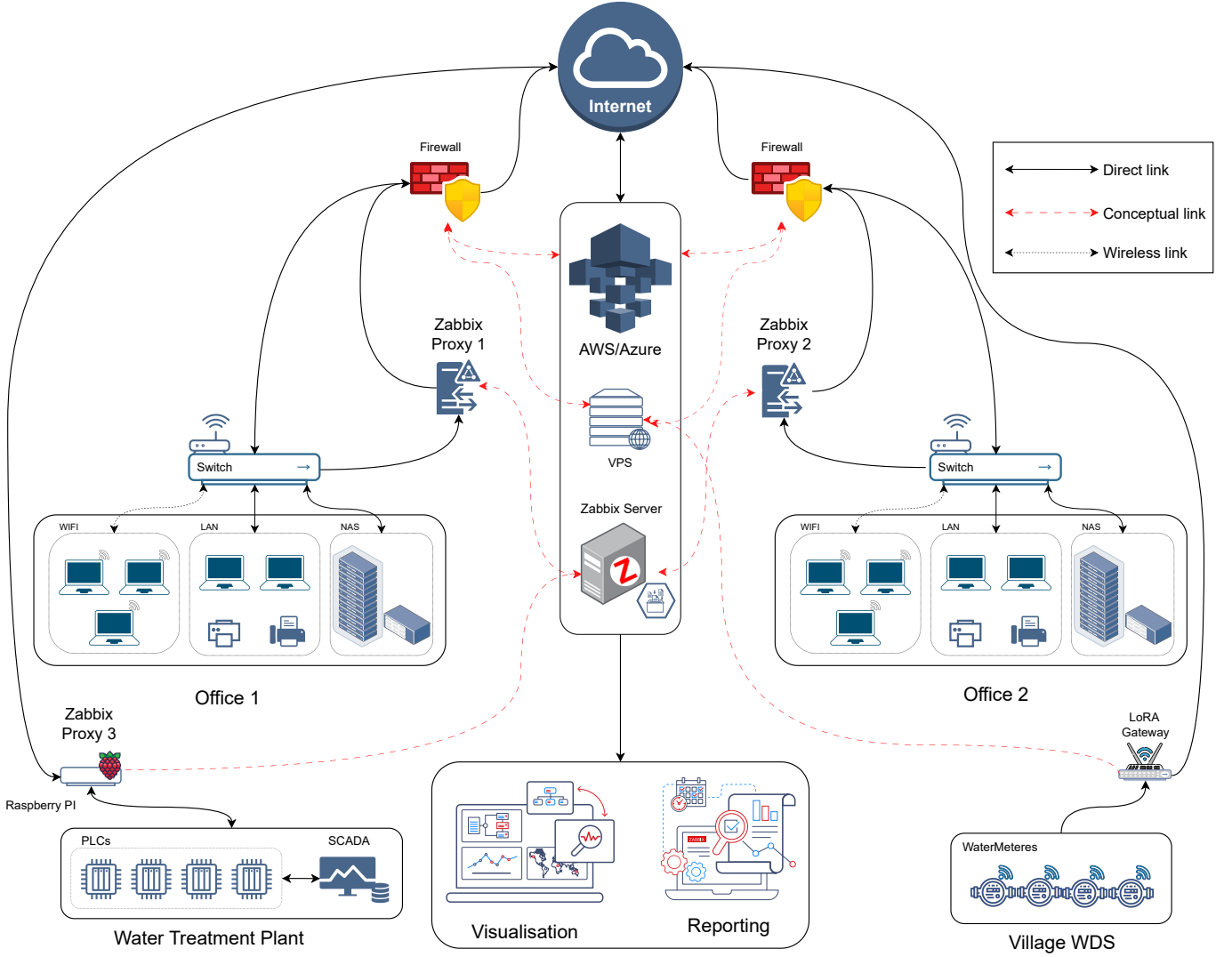
Fig. 3. VCD architecture with Zabbix and ML-based IDS for DT-enabled SME water systems

wide range of attacks while remaining lightweight and suitable for resource-constrained environments.

### D. AI/ML-Based Intrusion Detection System

As part of the proposed framework, we developed a machine learning-based IDS to improve the detection of cyber threats in smart water networks. This IDS is trained on the OD-IDS2022 dataset, which provides 1,031,916 labeled samples [18]. Each sample contains 82 features representing flow-based network data, including IP addresses, port numbers, protocol types, packet lengths, time durations, and flag behaviors. These records include normal traffic and 29 attack types such as DoS, brute force, SQL injection, RCE, hijacking, and reconnaissance. To simplify classification and reduce overfitting, we grouped the 29 attack classes into seven general categories, listed in Table II. This grouping keeps the detection meaningful while making the machine learning models easier to train and evaluate.

TABLE II
7-GROUP ATTACK CATEGORIZATION

| Group | Includes |
|---|---|
| BENIGN | BENIGN |
| DOS | DoS Hulk, Slowhttptest, GoldenEye, Slowloris, DDoS-* |
| BRUTEFORCE | Bruteforce-Web, Bruteforce-XSS, FTP/SSH-Patator, Web Brute Force |
| INJECTION | SQL/LDAP/SIP Injection, Web SQL Injection |
| HIJACKING | MITM, Hijacking |
| RCE | RFI, Exploit, Cmd Injection, Upload, Backdoor |
| OTHER | Infiltration, Bot, PortScan, Web XSS |

We proposed and implemented five machine learning models as part of the IDS component. All models use the same preprocessing pipeline: label encoding, numerical feature extraction, mutual information for feature selection, data normalization, and oversampling with SMOTE to balance class distribution.

*1) Random Forest Classifier:* The Random Forest (RF) model builds many decision trees from random subsets of the training data. Each tree gives a class prediction, and the final result is selected by majority voting. This is expressed in

Equation 1.

$$\hat{y} = \text{mode}\,(h_1(x), h_2(x), \ldots, h_T(x)) \tag{1}$$

where $h_t(x)$ is the prediction of tree $t$, and $T$ is the total number of trees.

*2) Tuned LightGBM Classifier:* LightGBM is a gradient boosting algorithm that builds trees sequentially to minimize prediction errors. It grows trees leaf-wise and uses a loss function with regularization, as shown in Equation 2.

$$\mathcal{L}^{(t)} = \sum_{i=1}^{n} L(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \Omega(f_t) \tag{2}$$

Here, $L$ is the loss function, $\hat{y}_i^{(t-1)}$ is the previous prediction, $f_t$ is the new decision tree, and $\Omega$ is the regularization term.

*3) Improved Ensemble Model (v1):* This model combines three base classifiers—Random Forest, LightGBM, and a Multi-layer Perceptron (MLP)—into a soft voting ensemble. It averages the class probabilities from each model and selects the class with the highest average score, as shown in Equation 3.

$$\hat{y} = \arg\max_{c} \left( \frac{1}{M} \sum_{m=1}^{M} P_m(c \mid x) \right) \tag{3}$$

where $P_m(c \mid x)$ is the probability of class $c$ predicted by model $m$, and $M$ is the number of models.

*4) Weighted Ensemble with Feature Engineering:* This model improves ensemble voting by assigning custom weights to each classifier and using new engineered features like packet length ratios and size variations. The prediction formula with weights is given in Equation 4.

$$\hat{y} = \arg\max_{c} \left( \sum_{m=1}^{M} w_m \cdot P_m(c \mid x) \right) \tag{4}$$

where $w_m$ is the weight assigned to model $m$, and $\sum w_m = 1$.

*5) Improved Ensemble (v2):* The final and most optimized model uses the same weighted voting as in Equation 4, but with improved components. These include:

- A deeper MLP with 3 hidden layers (256, 128, 64) and ReLU activation
- A tuned LightGBM with max depth = 10, 64 leaves, and learning rate = 0.05
- A larger Random Forest with 150 trees and class-balanced weighting

The ensemble weights are selected based on validation scores to ensure balanced detection across all classes, especially minority attacks like RCE and Hijacking.

*Note:* To improve model transparency, we use SHAP (SHapley Additive exPlanations), a method from cooperative game theory that attributes prediction changes to individual features. The SHAP value for a feature $i$ is calculated using:

$$\phi_i = \sum_{S \subseteq F \setminus \{i\}} \frac{|S|!(|F| - |S| - 1)!}{|F|!} \left[ f(S \cup \{i\}) - f(S) \right] \tag{5}$$

Here, $F$ is the full feature set, $S$ is a subset of features excluding $i$, and $f$ is the prediction function. SHAP values explain how much each feature contributes to the final prediction, helping operators understand the decision process of the IDS.

In summary, the AI-based IDS module strengthens the proposed framework by enabling real-time detection of various cyber threats using interpretable and resource-efficient machine learning models. In the following section, we present the experimental evaluation and performance results of the IDS models, along with the integration of Zabbix for continuous monitoring in the digital twin environment.

## IV. EXPERIMENTAL EVALUATION AND MODEL PERFORMANCE

This section presents the experimental evaluation of the proposed VCD for water systems. The validation includes two parts: real-time monitoring results using Zabbix, and the performance of machine learning models for intrusion detection.

### A. Monitoring Setup and Attack Simulation

The VCD was tested in a hybrid digital twin setup. Zabbix server was installed on a VPS, and several Raspberry Pi devices were installed in field locations like water plants and offices. These Raspberry Pis worked as Zabbix proxies and collected logs from IoT meters, PLCs, and SCADA systems.

To test the system, three types of cyberattacks were performed:

- *Nmap Scan (Reconnaissance)*: A stealth scan was launched using Nmap to find open ports. Suricata detected this scan and sent alerts to Zabbix. The traffic pattern showed abnormal packet behavior (Figure 4).
- *Brute Force (Hydra + SSH)*: An SSH brute-force attack was simulated using Hydra. Zabbix recorded many failed logins and increased CPU usage. Suricata also detected frequent access to port 22. Fail2Ban blocked the attacker's IP after too many failed attempts, as shown in Figures( 5, 6).
- *DoS (hping3)*: A SYN flood attack was done using hping3. It caused high CPU and memory usage. Zabbix showed this unusual behavior and created alerts, even when the logs were not clear.

These tests showed that the system can detect and respond to real cyberattacks using simple and open-source tools.

### B. Monitoring Indicators

Several indicators were collected from Zabbix to check the system behavior:

- *CPU and Memory Usage*: These increased during DoS attack and helped to detect it (Figure 7).

| Timestamp | Name | Value |
|---|---|---|
| 2025-03-14 10:45:02 | fast.log | 03/14/2025-10:23:01.719578 [**] [1:3400002:2] POSSBL PORT SCAN (NMAP -sS) [**] [Classification: Attempted Information Leak] [Priority: 2] {TCP} 192.168.88.242:61861 -> 192.168.88.254:22641 |
| 2025-03-14 10:45:02 | fast.log | 03/14/2025-10:23:01.718516 [**] [1:3400002:2] POSSBL PORT SCAN (NMAP -sS) [**] [Classification: Attempted Information Leak] [Priority: 2] {TCP} 192.168.88.242:61861 -> 192.168.88.254:33406 |
| 2025-03-14 10:45:02 | fast.log | 03/14/2025-10:23:01.714599 [**] [1:3400002:2] POSSBL PORT SCAN (NMAP -sS) [**] [Classification: Attempted Information Leak] [Priority: 2] {TCP} 192.168.88.242:61861 -> 192.168.88.254:1991 |
| 2025-03-14 10:45:02 | fast.log | 03/14/2025-10:23:01.713213 [**] [1:3400002:2] POSSBL PORT SCAN (NMAP -sS) [**] [Classification: Attempted Information Leak] [Priority: 2] {TCP} 192.168.88.242:61861 -> 192.168.88.254:25382 |
| 2025-03-14 10:45:02 | fast.log | 03/14/2025-10:23:01.711525 [**] [1:3400002:2] POSSBL PORT SCAN (NMAP -sS) [**] [Classification: Attempted Information Leak] [Priority: 2] {TCP} 192.168.88.242:61861 -> 192.168.88.254:2703 |
| 2025-03-14 10:45:02 | fast.log | 03/14/2025-10:23:01.711155 [**] [1:3400002:2] POSSBL PORT SCAN (NMAP -sS) [**] [Classification: Attempted Information Leak] [Priority: 2] {TCP} 192.168.88.242:61861 -> 192.168.88.254:50570 |

Fig. 4. logging attempt to the servers

| Name | Value |
|---|---|
| fail2ban.log | 2025-03-14 13:35:34,010 fail2ban.actions [4542]: NOTICE [sshd] Ban 192.168.88.1 |
| fail2ban.log | 2025-03-14 13:27:45,376 fail2ban.actions [4542]: NOTICE [sshd] Ban 192.168.88.242 |
| fail2ban.log | 2025-03-14 13:26:17,934 fail2ban.actions [3253]: NOTICE [sshd] Flush ticket(s) with iptables- |
| fail2ban.log | 2025-03-14 12:48:49,397 fail2ban.actions [3253]: NOTICE [sshd] Unban 192.168.88.242 |
| fail2ban.log | 2025-03-14 12:40:08,214 fail2ban.actions [3253]: WARNING [sshd] 192.168.88.242 already banned |
| fail2ban.log | 2025-03-14 12:40:08,212 fail2ban.actions [3253]: WARNING [sshd] 192.168.88.242 already banned |
| fail2ban.log | 2025-03-14 12:40:08,210 fail2ban.actions [3253]: WARNING [sshd] 192.168.88.242 already banned |
| fail2ban.log | 2025-03-14 12:40:08,209 fail2ban.actions [3253]: WARNING [sshd] 192.168.88.242 already banned |
| fail2ban.log | 2025-03-14 12:40:08,208 fail2ban.actions [3253]: WARNING [sshd] 192.168.88.242 already banned |
| fail2ban.log | 2025-03-14 12:40:08,207 fail2ban.actions [3253]: WARNING [sshd] 192.168.88.242 already banned |
| fail2ban.log | 2025-03-14 12:40:08,204 fail2ban.actions [3253]: WARNING [sshd] 192.168.88.242 already banned |

Fig. 5. Fail2Ban logs showing IP bans triggered by repeated failed SSH login attempts

- *Network Traffic (Upload/Download)*: Abnormal traffic helped detect Nmap and DoS attacks (Figure 8).
- *Dropped and Malformed Packets*: These increased during the flood attack.
- *Failed Login Attempts*: Zabbix tracked this for brute force detection, and Fail2Ban blocked the IP.
- *Suricata Alerts*: Number of alerts helped show which attack was happening.
- *Alert Time*: The system created alerts in a few seconds after the attack started.

### C. Machine Learning IDS Evaluation

Besides traditional detection, five machine learning models were tested using the OD-IDS2022 dataset. The goal was to

| Timestamp | Name | Value |
|---|---|---|
| 2025-03-14 12:39:02 | fast.log | 03/14/2025-12:38:43.792484 [**] [1:1000002:1] SSH Brute Force Attempt [**] [Classification: Attempted Administrator Privilege Gain] [Priority: 2] {TCP} 192.168.88.242:48658 -> 192.168.88.254:22 |
| 2025-03-14 12:38:02 | fast.log | 03/14/2025-12:37:32.925070 [**] [1:1000002:1] SSH Brute Force Attempt [**] [Classification: Attempted Administrator Privilege Gain] [Priority: 2] {TCP} 192.168.88.242:35306 -> 192.168.88.254:22 |
| 2025-03-14 12:37:02 | fast.log | 03/14/2025-12:36:32.965751 [**] [1:1000002:1] SSH Brute Force Attempt [**] [Classification: Attempted Administrator Privilege Gain] [Priority: 2] {TCP} 192.168.88.242:45532 -> 192.168.88.254:22 |

Fig. 6. Suricata alerts for SSH brute-force attack attempts showing repeated unauthorized access to port 22
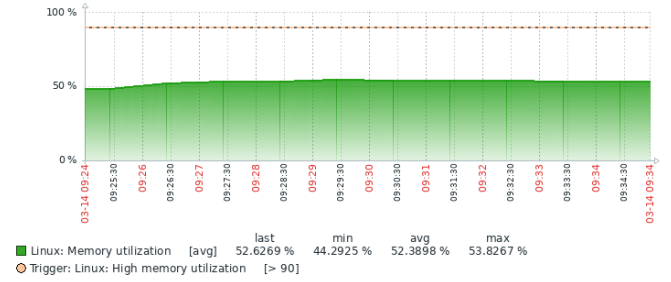


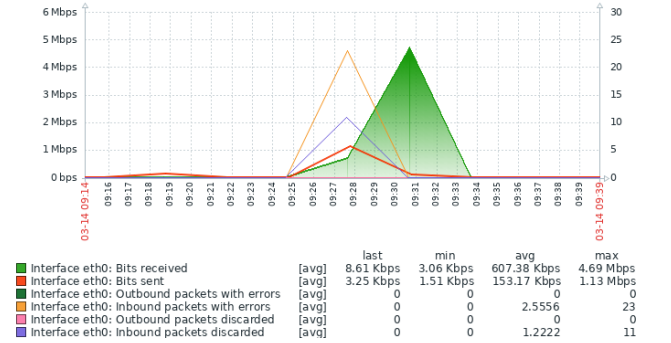Fig. 7. Memory usage monitoring under DDoS Attack



Fig. 8. Network monitoring under DDoS Attack

classify seven types of network traffic, including attacks like RCE, hijacking, and injection.

Table III shows the performance of each model. The best model was the improved ensemble (v2), which used Light-GBM, Random Forest, and MLP together.

TABLE III
COMPARISON OF IDS MODELS FOR 7-CLASS CATEGORIZATION
(OD-IDS2022 DATASET)

| Model | Acc. | Macro F1 | RCE F1 | HIJACK Rec. | Explainable | Gran. |
|---|---|---|---|---|---|---|
| Random Forest | 77.0% | 0.47 | 0.37 | 0.54 | ✓ SHAP | 30+ |
| LightGBM (Tuned) | 82.2% | 0.714 | 0.526 | 0.609 | – (addable) | 7 |
| Improved Ens. (v1) | 80.4% | 0.6645 | 0.55 | 0.73 | ✓ SHAP | 7 |
| Weighted Ens. + FE | 80.2% | 0.66 | 0.54 | 0.76 | ✓ SHAP | 7 |
| Improved Ens. (v2) | 92.0% | 0.88 | 0.86 | 0.87 | ✓ SHAP | 7 |

Table IV shows the full report for the best model. It gives good results in all classes, including small ones like injection. Figure 9 shows the confusion matrix.

TABLE IV
ENSEMBLE MODEL CLASSIFICATION REPORT

| Class | Precision | Recall | F1-score | Support |
|---|---|---|---|---|
| BENIGN | 0.91 | 0.93 | 0.92 | 2024 |
| BRUTEFORCE | 0.87 | 0.85 | 0.86 | 2377 |
| DOS | 0.96 | 0.97 | 0.96 | 6463 |
| HIJACKING | 0.84 | 0.87 | 0.85 | 2475 |
| INJECTION | 0.82 | 0.78 | 0.80 | 220 |
| OTHER | 0.92 | 0.93 | 0.93 | 14629 |
| RCE | 0.85 | 0.88 | 0.86 | 1812 |
| Accuracy | | | 0.92 | 30000 |
| Macro Avg | 0.88 | 0.89 | 0.88 | 30000 |
| Weighted Avg | 0.92 | 0.92 | 0.92 | 30000 |

This experiment confirms that the proposed VCD can detect and respond to cyberattacks in real time using both rule-based
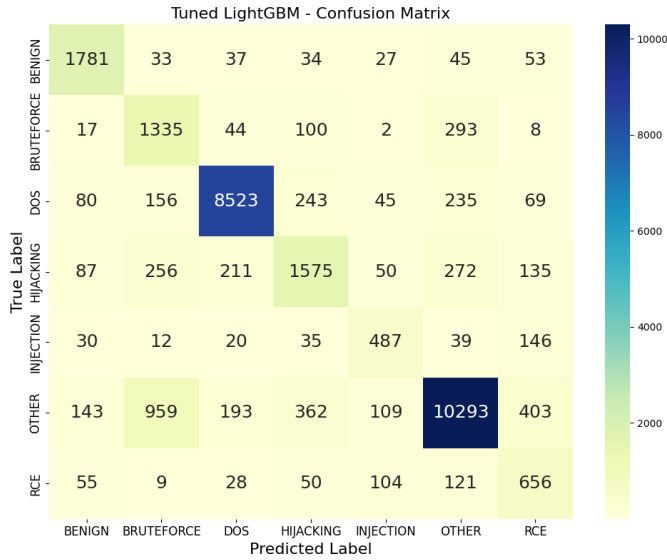
Fig. 9. Confusion matrix showing class-wise prediction performance across 7 traffic categories.

and AI-based tools. It works well even in small water systems with low-cost hardware.

## V. Conclusion and Future Work

This study proposed a VCD to improve the cybersecurity of DTs used in water networks, with a focus on SMEs. The solution combines free tools: Zabbix for live monitoring, Suricata as an IDS, and Fail2Ban to block repeated login attempts. Zabbix proxies were installed on Raspberry Pi units to collect data from SCADA, PLCs, and IoT sensors. We tested the system with simulated attacks (port scanning, brute-force on SSH, and DoS), and it responded correctly with alerts, log collection, and IP blocking.

The IDS part was developed using the OD-IDS2022 dataset (over one million records with 29 classes). We simplified the task by grouping the classes into 7 attack types. We tested five ML models, and the final version (v2) used a combination of RF, LGBM, and MLP with SHAP for explainability. This model gave the best results for detecting attacks like RCE, hijacking, and injection. The full framework is low-cost, supports real-time detection, and works well for small organizations without advanced computing systems.

For future work, we will explore LLMs to improve detection accuracy, reduce false positives, and classify threats more precisely. We also plan to integrate blockchain to protect data integrity and support trusted operations. These upgrades aim to create smarter and more secure water management systems.

## Acknowledgment

## References

[1] M. Homaei, O. Mogollón-Gutiérrez, J. C. Sancho, M. Ávila, and A. Caro, "A review of digital twins and their application in cybersecurity based on artificial intelligence," Artificial Intelligence Review, vol. 57, no. 8, jul 2024.

[2] M. Homaei, A. C. Lindo, J. C. S. N. nez, O. M. Gutiérrez, and J. A. Díaz, "The role of artificial intelligence in digital twin's cybersecurity," in XVII Reunión Española Sobre Criptología y Seguridad de La Información (RECSI), vol. 265, 2022, p. 133.

[3] A. Abbasi, F. Zaidi, and O. F. Rana, "A survey on cybersecurity in critical infrastructures: Approaches, challenges, and future directions," ACM Computing Surveys, vol. 54, no. 7, pp. 1–36, 2021. doi: 10.1145/3453158.

[4] J. Smith and A. Brown, "A review of cyber threats and security solutions for water distribution networks," Journal of Industrial Information Integration, vol. 15, pp. 100–110, 2019. doi: 10.1016/j.jii.2019.07.002.

[5] S. Yu and L. Liu, "Survey on cyber-physical system security in water sector: Threats and defense strategies," IEEE Access, vol. 9, pp. 78765–78779, 2021. doi: 10.1109/ACCESS.2021.3083795.

[6] A. Brown, C. Wu, and D. Wilson, "Anomaly detection in critical infrastructures: A survey of methods and challenges," IEEE Communications Surveys & Tutorials, vol. 24, no. 3, pp. 177–198, 2022. doi: 10.1109/COMST.2022.3141234.

[7] F. Tao and M. Zhang, "Digital twin shop-floor: A new shop-floor paradigm towards smart manufacturing," IEEE Access, vol. 5, pp. 20418–20427, 2018. doi: 10.1109/ACCESS.2017.2756069.

[8] A. Fuller, Z. Fan, C. Day, and C. Barlow, "Digital twin: Enabling technologies, challenges and open research," IEEE Access, vol. 8, pp. 108952–108971, 2020. doi: 10.1109/ACCESS.2020.2998358.

[9] A. Mirchi and K. Madani, "Water resources cyber-physical security: A review and future research directions," Water, vol. 12, no. 6, p. 1602, 2020. doi: 10.3390/w12061602.

[10] V. Kandasamy and M. Shankaran, "Survey on open-source IDS and monitoring solutions for critical infrastructures," International Journal of Network Security, vol. 23, no. 1, pp. 45–58, 2021. doi: 10.6633/IJNS.202101_23(1).06.

[11] S. Khan, H. Li, and R. Ahmad, "A review of predictive analytics in critical infrastructure management for sustainability," Sustainable Cities and Society, vol. 80, p. 103789, 2022. doi: 10.1016/j.scs.2022.103789.

[12] Z. Zhang, X. Chen, G. Zhao, and W. Gao, "A cyber-physical attack detection framework for digital twin-based water distribution systems," IEEE Internet of Things Journal, vol. 8, no. 9, pp. 7650–7660, 2021.

[13] D. Liu, M. Zhong, Y. Fu, and Y. Li, "Cybersecurity for digital twin-based smart water management systems," Journal of Water Resources Planning and Management, vol. 148, no. 4, p. 04022005, 2022. doi: 10.1061/(ASCE)WR.1943-5452.0001510.

[14] S. Qi, K. Wang, and H. Zhuang, "Cybersecurity risk assessment in digital twin-enabled smart water networks," Sensors, vol. 22, no. 18, p. 6905, 2022. doi: 10.3390/s22186905.

[15] S. Kumar, Y. Wu, and J. Li, "Mitigating cyber-attacks in digital twin environments of water distribution networks using anomaly detection and blockchain," IEEE Transactions on Industrial Informatics, 2023, Early Access.

[16] Y. Lin, R. Deng, and H. Chen, "Securing cyber-physical water distribution infrastructures: A digital twin-based intrusion detection approach," Journal of Network and Computer Applications, vol. 225, p. 103525, 2023. doi: 10.1016/j.jnca.2022.103525.

[17] M. Homaei, A. J. Di Bartolo, M. Ávila, O. Mogollón-Gutiérrez, and A. Caro, "Digital transformation in the water distribution system based on the digital twins concept," 2024.

[18] N. D. Patel, B. M. Mehtre, and R. Wankar, "Od-ids2022: generating a new offensive defensive intrusion detection dataset for machine learning-based attack classification," International Journal of Information Technology, vol. 15, no. 8, pp. 4349–4363, Sep. 2023. doi: 10.1007/s41870-023-01464-8.