

The Dark Side of Digital Twins: Adversarial Attacks on AI-Driven Water Forecasting

Mohammadhossein Homaei
Grupo de Ingeniería de Medios
Universidad de Extremadura
Mhomaein@alumnos.unex.es

Víctor González Morales,
Óscar Mogollón-Gutiérrez
Grupo de Ingeniería de Medios
Universidad de Extremadura
{victorgomo, oscarmg}@unex.es

Andrés Caro
Grupo de Ingeniería de Medios
Universidad de Extremadura
andresc@unex.es

Abstract—Digital twins (DTs) are improving water distribution systems by using real-time data, analytics, and prediction models to optimize operations. This paper presents a DT platform designed for a Spanish water supply network, utilizing Long Short-Term Memory (LSTM) networks to predict water consumption. However, machine learning models are vulnerable to adversarial attacks, such as the Fast Gradient Sign Method (FGSM) and Projected Gradient Descent (PGD). These attacks manipulate critical model parameters, injecting subtle distortions that degrade forecasting accuracy. To further exploit these vulnerabilities, we introduce a Learning Automata (LA) and Random LA-based approach that dynamically adjusts perturbations, making adversarial attacks more difficult to detect. Experimental results show that this approach significantly impacts prediction reliability, causing the Mean Absolute Percentage Error (MAPE) to rise from 26% to over 35%. Moreover, adaptive attack strategies amplify this effect, highlighting cybersecurity risks in AI-driven DTs. These findings emphasize the urgent need for robust defenses, including adversarial training, anomaly detection, and secure data pipelines.

Index Terms—Digital Twins, Artificial Intelligence, Cybersecurity, Adversarial Machine Learning Attack

I. INTRODUCTION

Digital twin technology has emerged as a critical driver of digital transformation across various industries, playing a crucial role in improving the accuracy and efficiency of cyber-physical systems. One area where this technology has been widely adopted is in water distribution networks. DTs enable real-time and accurate simulations of physical systems, facilitating enhanced decision-making and optimization of operations. Through data-driven and automated processes, this technology helps industries increase operational efficiency and utilize resources more sustainably [1], [2].

However, despite their many benefits, DTs are exposed to significant security challenges due to their complex structure and continuous connection to the internet. One of the most critical security threats is data and model poisoning, which can significantly compromise the performance of digital twin systems and lead to erroneous outcomes. Data poisoning involves manipulating input data that feeds into the system, while model poisoning involves tampering with ML models [3], [4]. These threats pose serious risks, especially for systems that rely on ML models for prediction and optimization, as they can degrade the accuracy of forecasts and increase operational costs [5].

In this context, AML attacks, particularly FGSM, have demonstrated the ability to disrupt ML models [6]–[8], and

PGD [9] consists of the iterative application of FGSM. These attacks introduce small perturbations to input data, drastically reducing model accuracy and misleading prediction systems. Given that ML models such as LSTM networks are extensively used for water consumption forecasting in distribution networks, such attacks can significantly impact the efficiency and effectiveness of these systems [10], [11].

This paper focuses on addressing the security challenges faced by DTs in water distribution networks and proposes innovative solutions to mitigate the risks of data and model poisoning. A security layer has been developed to safeguard the system against cyberattacks like FGSM and data poisoning, ensuring system integrity. The evaluation results demonstrate that the proposed platform improves the system's accuracy and efficiency, reduces operational costs, and supports intelligent decision-making. These solutions are vital for ensuring the sustainability of water resources and advancing digital transformation in the water sector.

A. Objective of the Article

This paper focuses on cybersecurity risks in DTs for water distribution networks, specifically in water consumption forecasting using time series data. It examines AML attacks, such as FGSM and PGD, and their impact on LSTM-based prediction models. To make attacks more complex and harder to detect, the study explores the use of LA. Additionally, it proposes mitigation strategies to strengthen DTs against data and model poisoning, improving system reliability, reducing costs, and supporting better decision-making in digital water management.

B. Paper Structure

This paper is structured as follows: Section II reviews related work and the motivation behind addressing cybersecurity in WDS. Section III presents the proposed DT platform and forecasting models. Section IV analyzes the impact of FGSM-based AML attacks on LSTM models. Section V applies the FGSM attack and evaluates forecasting accuracy. Section VI introduces a Learning Automata-based FGSM attack to improve stealth. Section VII proposes a Random Learning Automata strategy to enhance unpredictability. Section VIII discusses DT vulnerabilities and mitigation strategies. Finally, Section IX concludes the paper and outlines future directions.

Table I
DT PROJECTS IN WATER INDUSTRY WITH AI/ML/DL AND THEIR VULNERABILITIES TO POISONING ATTACKS

Project	Data Used	Model/Algorithm	ML/AI/DL Techniques	Vulnerability	Possible Attack Vectors	Mitigation Strategies
Ciliberti et al. (2021)	Pressure, flow, and asset data	AI models for leak detection, optimization	ML/AI for DMA optimization	High	Data poisoning of asset management systems	Data integrity verification, encrypted data channels
Bonilla et al. (2022)	Real-time pressure and flow rate data	GCNs	AI/ML for hydraulic state estimation	High	Data poisoning of real-time pressure and flow data	Data validation techniques, secure data pipelines
Zekri et al. (2022)	IoT sensor data, asset operation data	Multi-Agent Systems, AI-driven agents	Multi-agent reinforcement learning	Moderate-High	Model poisoning by corrupting reward system	Robust reward functions, anomaly detection
Matheri et al. (2022)	Wastewater treatment sensor data	Cyber-Physical Systems (CPS), AI optimization models	AI/ML for predictive maintenance	High	Data poisoning via corrupted sensor data	Secure real-time data transmission, predictive anomaly detection
Ramos et al. (2022)	Leakage detection, water usage data	Optimization algorithms	AI-driven optimization for water management	Moderate	Tampered leakage data poisoning	Cryptographic methods for data validation
Henriksen et al. (2022)	Hydrological data, climate models	ML Models	ML for climate adaptation	Moderate	Tampered hydrological data affecting water management	Redundant data sources, data quality audits
Savic (2022)	Water usage patterns	AI anomaly detection models	AI/ML for anomaly detection	High	Model poisoning via corrupted anomaly detection data	AI model validation, adversarial training
Pedersen et al. (2022)	Water level sensors, drainage data	Hydraulic models with DTs	ML for error classification	Moderate	Data poisoning from water level sensors	Redundant sensor validation, secure data transmission
Valencia Smart Water (2023)	SCADA, customer feedback	Hydraulic models with real-time optimization	AI/ML for pressure management	Moderate	Tampered sensor inputs poisoning	Secure data transmission, blockchain verification
Sabesp Digital Twin (2023)	IoT sensor data, remote monitoring	ML for anomaly detection	AI/ML for fault detection	High	Data poisoning from faulty sensors or pump data manipulation	Real-time anomaly detection, secure IoT devices
Tarragona Water Consortium (2023)	Hydraulic models, real-time sensors	Live simulations for predictive maintenance	ML for predictive maintenance	Moderate	Data manipulation through compromised sensors	Redundant sensor data validation systems
Smart Water Grid in Gaula (2023)	Leakage detection, water usage data	Digital twin with real-time optimization	AI-driven optimization for water loss prevention	Moderate-High	Data poisoning from compromised sensors	Blockchain-based validation, real-time monitoring
Water Research Foundation AI/ML Project (2023)	Utility performance data	ML models for prediction	AI/ML for performance optimization	High	Model poisoning via contaminated datasets	Regular model audits, federated learning
Menapace et al. (2024)	Pressure sensor data	GNNs	DL/ML for pressure estimation	High	Data poisoning through compromised sensor data	Data cross-validation, hybrid training with anomaly detection

II. RELATED WORKS AND MOTIVATION

A. Related Works

The increasing reliance on AI and ML in water distribution networks has led to significant advancements in digital twin technology. However, these AI-driven systems also introduce cybersecurity vulnerabilities, particularly AML threats. Several studies have investigated the application of DTs, AI, and cybersecurity in water management.

One of the key areas of research involves integrating AI models, such as LSTM networks, into water forecasting systems. Studies like those of [10], [11] have demonstrated the effectiveness of LSTM models in accurately predicting water consumption patterns. However, these models remain susceptible to adversarial attacks that can degrade their predictive accuracy.

Another line of research focuses on securing AI-based water management systems from cyber threats. [7], [8] examined adversarial attacks, including the FGSM, on AI-driven infrastructure, demonstrating how even minor perturbations in input data can significantly impact model predictions. Such vulnerabilities highlight the need for robust cybersecurity strategies in AI-powered DTs.

Additionally, various digital twin projects in the water industry have explored advanced ML techniques for enhanced system monitoring. [12] integrated graph convolutional networks (GCNs) within DTs for hydraulic state estimation, improving real-time system analysis. Similarly, [13] applied graph neural networks (GNNs) for sensor placement optimization, highlighting the growing intersection of AI and DTs in water distribution networks.

Despite these advancements, limited research has focused on the impact of AML threats on digital twin systems in the water industry. Existing studies primarily address general cybersecurity concerns or specific ML vulnerabilities but do not comprehensively analyze how adversarial attacks can compromise water forecasting accuracy. This gap underscores the necessity of further research on securing AI-based DTs from adversarial threats.

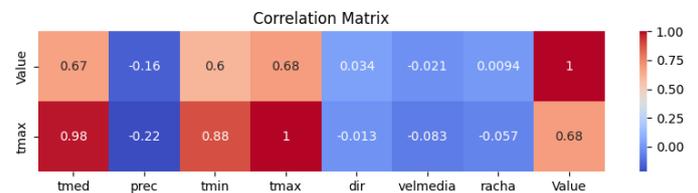


Figure 1. Correlation Matrix based on the parameters

B. Motivation for the Study

DT in water distribution networks enhances infrastructure management through predictive analytics, improving efficiency and sustainability. However, AI-driven forecasting models, particularly LSTMs, are highly susceptible to adversarial ML attacks.

This study addresses key challenges:

- **LSTM vulnerability:** Water consumption forecasting models, reliant on LSTMs, are prone to adversarial attacks like FGSM, degrading accuracy and increasing operational costs.
- **Insufficient security in DTs:** Most implementations lack robust cybersecurity measures to protect AI models from targeted attacks.
- **Expanding cybersecurity risks:** IoT-connected DTs enlarge the attack surface, making AI integrity crucial for system reliability.
- **Real-world impact:** Adversarial AI manipulation could disrupt water allocation, pressure control, and leak detection, compromising public utilities.

III. PROPOSED DT AND FORECASTING MODELS

CAUCCES is a DT platform designed to enhance water distribution through real-time monitoring, predictive analysis, and data-driven decision-making. It has been developed in collaboration with the Media Engineering Group at the University of Extremadura and Ambling Ingeniería y Servicios, S.L [14]. The platform integrates IoT sensors, AI-based forecasting, and secure data management to improve efficiency, minimize water loss, and ensure stable distribution.

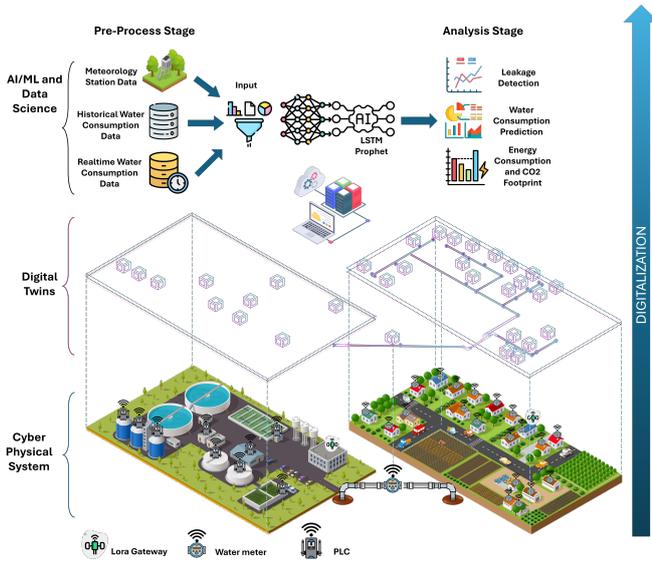


Figure 2. DT Platform in the Water Distribution Networks

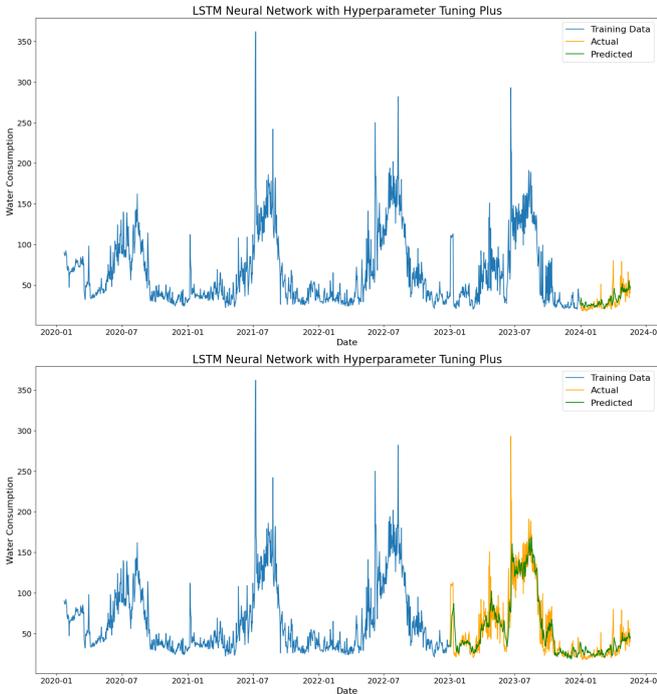


Figure 3. Water consumption forecasting via LSTM for 6(top) and 18 months(bottom)

Traditional water networks lack smart monitoring and forecasting, making them vulnerable to challenges such as aging infrastructure and environmental impacts. CAUCCES addresses these issues by continuously gathering data, utilizing reliable communication technologies, and optimizing scheduling for better maintenance and operation. This creates a real-time digital replica of the water system, enabling early detection and prevention of potential problems before they affect service.

A. Forecasting Results

Figure 3 presents the forecasting results of water consumption using the UV-LSTM models.

The following algorithm outlines the steps for training an LSTM model for water consumption prediction:

Algorithm 1 LSTM for Water Consumption Prediction

- 1 **Initialize parameters:**
- 2 Define the number of LSTM units (neurons), learning rate, and epochs
- 3 Initialize weight matrices W_f, W_i, W_C, W_o , and bias vectors b_f, b_i, b_C, b_o
- 4 **Preprocess input data:**
- 5 Normalize water consumption and meteorological data using Min-Max scaling
- 6 Divide the dataset into training, validation, and testing sets
- 7 Create sequences of input data X and target values Y
- 8 Reshape input data X to (num_samples, sequence_length, num_features)
- 9 **Model Training:**
- 10 Initialize cell state C_0 and hidden state h_0 to zeros
- 11 **for each epoch do**
- 12 **for each batch in the training data do**
- 13 **for each time step t in the input sequence do**
- 14 Compute Forget Gate: $f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$
- 15 Compute Input Gate: $i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$
- 16 Compute Candidate Cell State: $\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$
- 17 Update Cell State: $C_t = f_t \cdot C_{t-1} + i_t \cdot \tilde{C}_t$
- 18 Compute Output Gate: $o_t = \sigma(W_o \cdot [h_{t-1}, x_t] + b_o)$
- 19 Update Hidden State: $h_t = o_t \cdot \tanh(C_t)$
- 20 **end for**
- 21 Compute output (predicted water consumption): $y_{pred_t} = \text{Dense}(h_t)$
- 22 Calculate batch loss: $\text{MSE}(y_{pred_t}, y_{true_t})$
- 23 **Backpropagation through time (BPTT):**
- 24 Calculate gradients of Loss w.r.t weights and biases
- 25 Update W_f, W_i, W_C, W_o and b_f, b_i, b_C, b_o using an optimizer
- 26 **end for**
- 27 Evaluate model on the validation set after each epoch
- 28 **end for**
- 29 **Model Evaluation:**
- 30 Test the model on the testing set
- 31 Calculate and report performance metrics: RMSE and MAPE
- 32 **Model Deployment:**
- 33 Save the trained model for future use
- 34 Deploy the model for real-time or batch water consumption prediction

IV. AML ATTACKS ON FORECASTING

AML is an emerging field within the broader domain of ML, focusing on designing and evaluating models that are vulnerable to adversarial inputs. These inputs, also called adversarial examples, are carefully crafted perturbations designed to mislead the model into making incorrect predictions while appearing normal to humans. Such perturbations are typically small enough to go unnoticed in the input data but large enough to degrade the model's performance significantly. AML is crucial in various applications, particularly in domains where high reliability is required, such as healthcare, autonomous systems, and financial forecasting.

The *FGSM* is one of the foundational techniques in the adversarial attack literature. It is a white-box attack method, meaning the attacker has full access to the model, including its architecture and parameters. FGSM exploits the gradient of the loss function to the input features. Computing the gradient generates small perturbations to the input that maximally increase the model's loss, leading to erroneous predictions. The perturbation is scaled by a factor ϵ , which controls its magnitude. The mathematical expression for FGSM is given by:

$$X_{adv} = X + \epsilon \cdot \text{sign}(\nabla_X J(\theta, X, y)) \quad (1)$$

Where:

- X_{adv} is the adversarial input generated by the FGSM attack.
- X is the original input data (e.g., daily temperature and water consumption records).

- ϵ is the perturbation magnitude, determining the intensity of the attack.
- $J(\theta, X, y)$ represents the loss function, with θ denoting the model parameters and y being the true output label.
- $\nabla_X J(\theta, X, y)$ is the gradient of the loss function for the input data.

Meanwhile, the core formula for generating adversarial examples using PGD is as follows:

$$X_{\text{adv}}^{(t+1)} = \text{clip}_{[X-\epsilon, X+\epsilon]} \left(X_{\text{adv}}^{(t)} + \alpha \cdot \text{sign}(\nabla_{X_{\text{adv}}^t} J(X_{\text{adv}}^{(t)}, y)) \right) \quad (2)$$

Where:

- $X_{\text{adv}}^{(t)}$ refers to the generation of adversarial attacks at the time step t in the iterative sequence of the generation of FGSM attacks. At $t = 0$, it would be the input (daily water consumption or temperature).
- $X_{\text{adv}}^{(t+1)}$ represents the adversarially perturbed input at time step $t + 1$.
- ϵ is the perturbation factor determining the magnitude of the adversarial noise.
- $\nabla_{X_{\text{adv}}^t} J(\theta, X_{\text{adv}}^t, y)$ represents the gradient of the loss function concerning the input at time step t .
- $J(\theta, X_{\text{adv}}^t, y)$ is the loss function, and y is the true target value (water consumption in the following days).
- α is the step size at each iteration.
- $\text{sign}(\cdot)$ returns the sign of each component.
- The clip function projects the example into the ϵ -ball around x and enforces the valid data range.

By applying this formula, the adversarial input is slightly modified, causing the model to make incorrect predictions. This technique is both simple and effective, making it widely used in adversarial attack research.

V. APPLYING FGSM TO THE LSTM MODEL FOR WATER CONSUMPTION FORECASTING

The LSTM model, previously implemented to forecast water consumption based on daily water usage and temperature records, can be tested for robustness using FGSM-based adversarial attacks. The model's temporal nature, handling sequences of time-series data, makes it an interesting case for adversarial attacks, as small changes in the input sequence could propagate and lead to substantial forecasting errors over time.

To integrate FGSM into the LSTM framework, the first step involves computing the gradient of the loss function concerning the input sequence. Since LSTM models handle time series data, the inputs consist of daily historical records, including temperature and water consumption. The FGSM attack aims to perturb this input sequence so that the LSTM model's forecast deviates significantly from the actual future water consumption values.

The adversarial input X_{adv} is generated by adding a small perturbation ϵ to each feature (temperature and water consumption) at each time step in the input sequence. The sign of the gradient of the loss function determines the perturbation. After generating the adversarial input, the perturbed input sequence is fed into the LSTM model to assess how robust the model is against adversarially perturbed data.

Once the FGSM attack is applied, the impact on the LSTM model's predictive accuracy is evaluated. As shown in table II, even for small perturbation values ($\epsilon=0.001$), the Mean Absolute Error (MAE) and Root Mean Squared Error (RMSE) increase slightly, indicating a reduction in forecast precision. As ϵ increases, the error metrics grow significantly, with MAE exceeding 18 and MAPE surpassing 35% for ($\epsilon=0.01$). This demonstrates the model's vulnerability to small adversarial modifications in the input sequence.

Table II
RESULTS OF FGSM ATTACK WITH DIFFERENT ϵ VALUES ON LSTM MODELS

Model	ϵ	MAE	RMSE	MAPE
LSTM	0.0	12.334	20.653	24.419
LSTM	0.001	12.660	20.861	25.261
LSTM	0.005	13.959	21.718	28.609
LSTM	0.008	14.926	22.385	31.104
LSTM	0.01	15.567	22.840	32.754
LSTM+	0.0	12.329	20.449	25.032
LSTM+	0.001	12.758	20.726	26.108
LSTM+	0.005	14.462	21.879	30.378
LSTM+	0.008	15.721	22.783	33.531
LSTM+	0.01	16.551	23.400	35.606

To further evaluate the model's resilience, Projected Gradient Descent (PGD) is applied. Unlike FGSM, PGD refines the perturbation iteratively, leading to a stronger adversarial impact. Table III shows that for $\epsilon=0.01$, the MAE increases at a similar rate as FGSM, but for lower values of ϵ (e.g., 0.005 and 0.008), the prediction error already exhibits a steeper increase in RMSE and MAPE compared to FGSM. This suggests that even at intermediate perturbation levels, PGD induces more severe deviations in the LSTM model's forecasts.

Table III
RESULTS OF PGD ATTACK WITH DIFFERENT EPSILON VALUES ON LSTM MODELS

Model	ϵ	MAE	RMSE	MAPE
LSTM	0.0	12.334	20.653	24.419
LSTM	0.001	12.660	20.861	25.261
LSTM	0.005	13.960	21.720	28.615
LSTM	0.008	14.933	22.390	31.125
LSTM	0.01	15.579	22.848	32.790
LSTM+	0.0	12.329	20.449	25.032
LSTM+	0.001	12.759	20.726	26.108
LSTM+	0.005	14.467	21.884	30.395
LSTM+	0.008	15.741	22.797	33.588
LSTM+	0.01	16.589	23.424	35.708

Comparing both attacks, PGD consistently leads to higher errors at every tested epsilon value. While FGSM causes a steady degradation in model accuracy, PGD intensifies this effect by iteratively optimizing the perturbation, making it more effective in misleading the model. Notably, for $\epsilon=0.005$, the difference between FGSM and PGD is already evident in all error metrics, particularly in RMSE.

VI. LA-BASED UNDETECTABLE FGSM ATTACK

This section presents a learning automata-based approach for dynamically adjusting the perturbation size in the FGSM attack on LSTM models [15]. The primary objective is to improve attack stealth while maintaining its effectiveness in reducing forecasting accuracy.

A. ϵ Selection with Learning Automata

The learning automata mechanism selects an optimal ϵ value from a predefined set:

$$\epsilon \in \{0.0001, 0.0005, 0.001, 0.0025, 0.005\} \quad (3)$$

Each ϵ action has an associated probability, initialized equally and updated iteratively based on attack performance.

B. Probability Update Mechanism

The reward and penalty factors guide the probability updates:

- If the attack increases the MAPE within a controlled range, between 30 percent and 50 percent, the selected ϵ is rewarded.
- If MAPE exceeds 100 percent, the attack is considered too aggressive, and the probability of using that epsilon is penalized.

The probability update rule is given by:

$$P(a_t) = P(a_t) + r \cdot (1 - P(a_t)), \quad \text{if rewarded} \quad (4)$$

$$P(a_t) = P(a_t) \cdot (1 - p), \quad \text{if penalized} \quad (5)$$

where r is the reward factor, and p is the penalty factor.

C. Delayed Input Strategy

A delayed poisoning strategy introduces an artificial delay in using adversarial examples. Instead of applying the perturbation immediately, the adversarial inputs from previous iterations are stored and used after a fixed delay. This gradual attack approach helps to avoid abrupt changes, making the attack harder to detect.

D. Experimental Results

Using learning automata, the ϵ values were adjusted iteratively to maintain an effective yet undetectable attack. The MAPE progression over iterations showed a smooth increase, avoiding sharp variations, as shown in 4. The probability evolution of different ϵ values demonstrated the learning automata's ability to converge toward optimal attack parameters.

Figure 5 illustrates how an FGSM attack based on Learning Automata (LA) can remain imperceptible to a human observer. Unlike conventional perturbations that follow a monotonic increase in magnitude, the penalty and reward mechanism in LA introduces variability, preventing a straightforward detection pattern. This alternation in perturbation intensity adds a layer of randomness that disrupts the usual correlation between distortion and detectability. As a result, the attack does not exhibit a consistently increasing trend, making it more challenging to distinguish from natural fluctuations in the data. This adaptive behavior enhances the stealthiness of the adversarial perturbations, posing greater difficulties for both manual and automated detection methods.

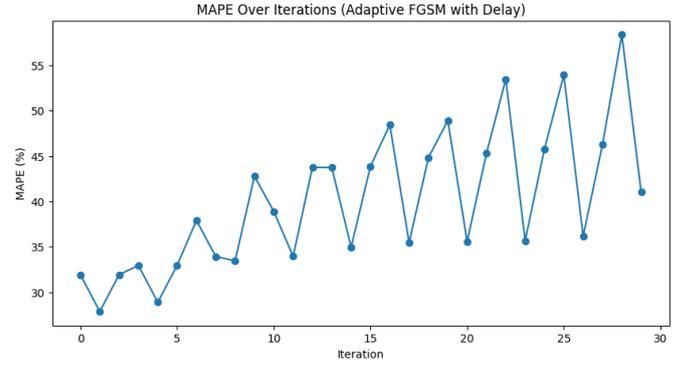


Figure 4. Fluctuation of the epsilon variable along the iterations in an LA-based FGSM Attack

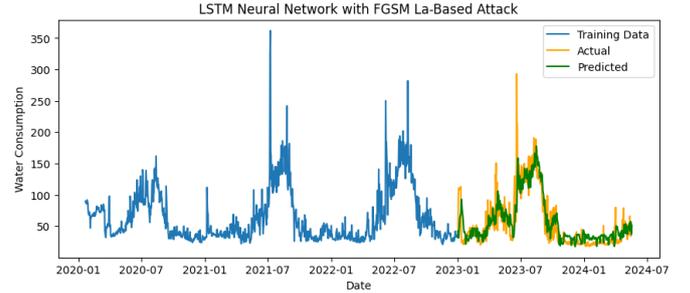


Figure 5. Hidden LA-based FGSM Attack

VII. RANDOM LEARNING AUTOMATA FOR FGSM ATTACK OPTIMIZATION

This section presents a Random Learning Automata (RLA)-based approach for dynamically adjusting the perturbation size in the FGSM attack on LSTM models. The key idea is to select multiple ϵ values per iteration instead of a single value, which helps improve attack stealth and avoid detection by the forecasting model [15].

A. ϵ Selection Using Random Learning Automata

In contrast to standard learning automata, RLA selects a random combination of ϵ values from a predefined set instead of a single value:

$$\epsilon \in \{0.0001, 0.0005, 0.001, 0.0025, 0.005\} \quad (6)$$

At each iteration, a subset of ϵ values is selected with probabilities determined by:

$$P(a_t) = \{P_1, P_2, \dots, P_n\}, \quad \sum_{i=1}^n P_i = 1 \quad (7)$$

where P_i represents the probability of selecting $\epsilon \in \epsilon_i$.

B. Probability Update Mechanism

To ensure the adaptive selection of ϵ values, the probability update mechanism follows:

$$P(a_t) = P(a_t) + r \cdot (1 - P(a_t)), \quad \text{if rewarded} \quad (8)$$

$$P(a_t) = P(a_t) \cdot (1 - p), \quad \text{if penalized} \quad (9)$$

where:

- r is the reward factor, controlling how quickly successful epsilon values gain priority.
- p is the penalty factor, reducing the probability of unsuccessful ϵ values.

The probabilities are normalized after every update:

$$P(a_t) = \frac{P(a_t)}{\sum_{j=1}^n P_j} \quad (10)$$

C. Multi ϵ Selection Strategy

Instead of choosing only one ϵ per iteration, RLA selects two or three ϵ values and applies them simultaneously:

$$\epsilon_{\text{chosen}} = \{\epsilon_i, \epsilon_j\}, \quad \text{where } i, j \in \{1, 2, 3, 4, 5\} \text{ and } i \neq j. \quad (11)$$

The number of selected ϵ values varies at each iteration and follows:

$$k \sim \mathcal{U}\{1, 3\} \quad (12)$$

where $\mathcal{U}(1, 3)$ represents a uniform distribution selecting either 1 or 2 ϵ values per iteration.

D. MAPE-Based Reward and Penalty System

The attack's effectiveness is measured using the MAPE. The MAPE is calculated as:

$$\text{MAPE} = \frac{100}{n} \sum_{t=1}^n \left| \frac{y_t - \hat{y}_t}{y_t} \right|, \quad y_t \neq 0 \text{ for all } t. \quad (13)$$

where:

- y_t is the actual water consumption value at time t .
- \hat{y}_t is the adversarially perturbed prediction.
- n is the total number of test samples.

The probability update mechanism follows these rules:

- If the attack increases MAPE within the range $30\% < \text{MAPE} < 50\%$, the selected ϵ values are rewarded.
- If MAPE exceeds 100%, the attack is too aggressive, leading to a penalty.
- If the MAPE increase per iteration is too high (above a threshold), a moderate penalty is applied.

To prevent abrupt changes in the attack, the penalty factor is adjusted dynamically:

$$p_{\text{adaptive}} = \begin{cases} 3p, & \text{if } \text{MAPE} > 100\% \\ 1.5p, & \text{if } \Delta\text{MAPE} > 5\% \\ p, & \text{otherwise} \end{cases} \quad (14)$$

E. Delay-Based Adversarial Example Storage

To further enhance stealth, adversarial examples are stored and applied after a delay. Instead of using the perturbed input immediately, RLA waits for a iterations before introducing the modified input into the model:

$$X_{\text{input}}^{(t)} = X_{\text{adv}}^{(t-a)} \quad (15)$$

where $X_{\text{adv}}^{(t-a)}$ is the adversarially generated input from a iterations ago. This delayed input poisoning strategy prevents sudden changes, making the attack harder to detect.

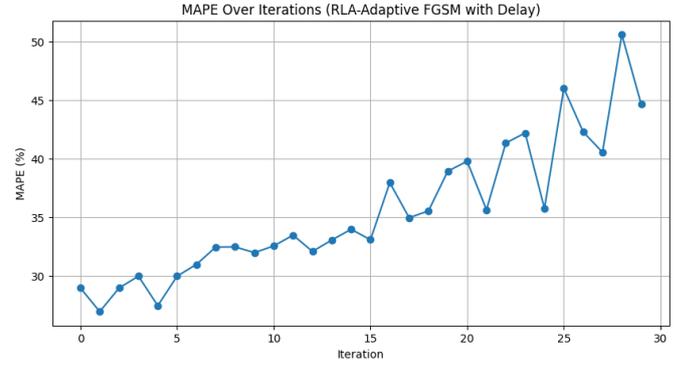


Figure 6. Fluctuation of the epsilon variable along the iterations in a RLA-based FGSM Attack

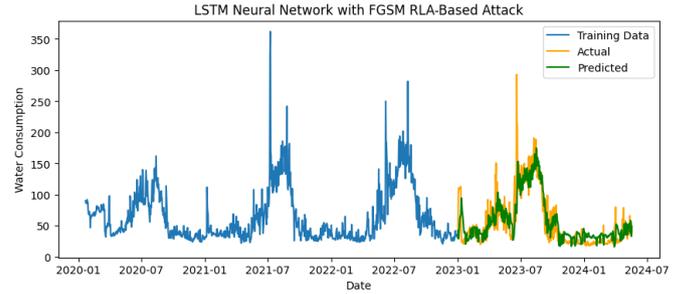


Figure 7. Hidden RLA-based FGSM Attack

F. Experimental Results

The FGSM attack based on Random Learning Automata (RLA) introduces a higher degree of unpredictability in the perturbation process. Unlike structured learning mechanisms, RLA selects ϵ values in a stochastic manner, influenced by random penalty and reward adjustments, as shown in 6. This randomness disrupts the formation of any discernible pattern, making the attack appear chaotic in nature. As a result, adversarial perturbations exhibit greater variation across different instances, reducing the likelihood of detection through conventional anomaly-based methods.

The key characteristic of RLA is its reliance on randomness rather than deterministic adaptation. The ϵ values fluctuate in a manner reminiscent of natural chaotic systems, where no two perturbations follow an exact progression. This irregularity prevents straightforward pattern recognition, complicating both manual and automated defense mechanisms. By embracing randomness as a core feature, the attack achieves a higher level of stealth, leveraging the unpredictability inherent in its learning process to bypass detection frameworks.

The impact of this randomness can be observed in the graphical representation of the ϵ values throughout the attack process. This variability makes it challenging to establish a clear boundary between adversarial and legitimate samples. Figure 7 illustrates this phenomenon, highlighting how the fluctuating nature of epsilon contributes to the stealthiness of the attack.

VIII. MITIGATION STRATEGIES

Protecting digital twin (DT) systems from adversarial attacks requires a combination of cybersecurity measures, data

Table IV
MITIGATION STRATEGIES FOR ADVERSARIAL ATTACKS IN DIGITAL TWIN FORECASTING MODELS

Aspect	Data / System	Technique / Method	Potential Vulnerability	Risk Level	Possible Attack Vectors	Mitigation Strategies
LoRa Encryption	Meter data, device provisioning	Built-in AES-128	Poor key management, reuse of keys	Moderate	Key guessing, eavesdropping on packets	Frequent key rotation, secure key distribution, proper implementation of AES-128
Meter-to-Gateway Sec	LoRa meter to LoRa gateway link	Mutual authentication, secure join procedures	Replay or spoofing of meter credentials	High	Impersonation of valid meters, data injection	Challenge-response protocols, nonce usage, short-lived session keys
Firmware Integrity	On-device firmware (water meters, gateways)	Signed or hashed firmware updates	Unauthorized firmware modifications	High	Malicious updates, remote code execution	Secure OTA updates with signature checks; regular patching
Gateway-to-Server	Data in transit from gateway to server	TLS/SSL, VPN, or private lines	Man-in-the-middle attacks	Moderate	Traffic interception, unauthorized data reading	Encrypted communication tunnels, certificate-based authentication
Net Monitoring & IDS	LoRa gateway and backend network	Intrusion Detection / Prevention Systems	Undetected brute-force, scanning attempts	Moderate	Malicious traffic patterns, repeated authentication failures	Automated anomaly detection, real-time threat response
ChirpStack Security	ChirpStack network server, application server	Role-based access, secure APIs	Misconfiguration, weak API keys	Moderate	Unauthorized device provisioning, data leakage	Secure API endpoints, regular security audits, minimal privilege policies
Database Security (PostgreSQL)	Stored water consumption records	Encryption at rest, access controls	Unauthorized DB access or tampering	High	SQL injection, stolen credentials	Strict role management, periodic audits, row-level security
End-to-End Encryption	Full data flow from meter to final storage	Consistent encryption in transit and at rest	Partial encryption gaps	Moderate	Plaintext exposure at intermediate hops, data sniffing	Holistic encryption approach, verifying encryption at all layers
AML Training for AI	Forecasting model (e.g., LSTM)	Incorporating FGSM/PGD examples into training	Model easily fooled by small perturbations	High	Data poisoning, gradient-based adversarial attacks	Model retraining on adversarial samples, ensemble methods
Domain-based Constraints	Forecasts	Physical/hydraulic plausibility checks	Acceptance of impossible meter readings	Moderate	Sudden large outliers can corrupt predictions	Filter or flag data outside valid usage/pressure thresholds
Real-time Anomaly Detection	Incoming meter data stream	Isolation Forest, One-Class SVM	Persistent adversarial or sensor tampering	High	Silent data drift, gradual poisoning	Trigger alerts when usage deviates from historical/seasonal norms
Gradient Masking & Model Randomization	Neural network layers	Adding noise, dropout, random inference steps	Straightforward gradient-based attack	Moderate	White-box adversary calculates precise gradients	Stochastic layers obscure gradients, raising attack complexity
Key Management & Regular Audits	All cryptographic operations	Rotating keys, HSM usage, compliance checks	Stolen or expired keys, unpatched systems	Moderate	Privilege escalation, extended infiltration	Automated key rotation, routine compliance (ISO/IEC 27001), robust backup/redundancy

integrity techniques, and machine learning defenses (Table IV). Key strategies focus on strengthening AI models against manipulation, ensuring secure data transmission, and implementing real-time monitoring. AI-based anomaly detection can identify suspicious activity, while adversarial training improves model robustness. Secure encryption, authentication protocols, and strict access controls help prevent unauthorized access and data tampering. Additionally, continuous auditing and compliance with cybersecurity standards ensure long-term resilience. By integrating these strategies, DTs can maintain reliable forecasting and decision-making even in the presence of cyber threats.

ACKNOWLEDGEMENT

This initiative is carried out within the framework of the funds from the Recovery, Transformation, and Resilience Plan, financed by the European Union (Next Generation) – National Institute of Cybersecurity (INCIBE), as part of project C107/23: "Artificial Intelligence Applied to Cybersecurity in Critical Water and Sanitation Infrastructures."

IX. CONCLUSION AND FUTURE WORK

This study shows that although AI-based DTs are helpful for water forecasting and resource management, they can still be attacked. Small, hidden changes in the data—called adversarial attacks—can reduce LSTM accuracy, increase costs, and damage trust. One serious method is AML, which can poison the system quietly. Our research highlights this hidden danger, which is important for city infrastructure and water systems where wrong decisions can have big impacts. Our method using learning automata can adapt to and follow the monthly and seasonal fluctuations in water consumption patterns.

Next, we plan to use Zabbix for real-time monitoring of things like sensor status, unusual forecasts, and network traffic. With smart alert settings, Zabbix can find problems fast and react automatically. In future work, we also want to use federated learning to avoid having one weak point and explore using multiple models together for stronger defense. These steps will help protect DTs and keep water systems safer from cyber threats.

REFERENCES

- [1] A.-J. Wang, H. Li, Z. He, Y. Tao, H. Wang, M. Yang, D. Savic, G. T. Daigger, and N. Ren, "Digital twins for wastewater treatment: A technical review," *Engineering*, vol. 36, p. 21–35, May 2024. [Online]. Available: <http://dx.doi.org/10.1016/j.eng.2024.04.012>
- [2] H. Beji and M. Lade, "Impact of digital transformation on carbon emissions reductions in the water industry," in *Lecture Notes in Energy*. Springer International Publishing, 2022, pp. 117–127.
- [3] M. Homaei, A. J. Di Bartolo, M. Ávila, Óscar Mogollón-Gutiérrez, and A. Caro, "Digital transformation in the water distribution system based on the digital twins concept," 2024. [Online]. Available: <https://arxiv.org/abs/2412.06694>
- [4] M. H. Homaei, A. C. Lindo, J. C. S. Núñez, O. M. Gutiérrez, and J. A. Díaz, Eds., *The role of Artificial Intelligence in Digital Twin's Cybersecurity*. Editorial Universidad de Cantabria, Sep. 2022.
- [5] M. Homaei, Óscar 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, p. 201, Jul. 2024. [Online]. Available: <https://doi.org/10.1007/s10462-024-10805-3>
- [6] E. Coda, B. Clymer, C. DeSmet, Y. Watkins, and M. Girard, "Universal fourier attack for time series," *IEEE Open Journal of Signal Processing*, vol. 5, p. 858–866, 2024.
- [7] S. U. Khan, M. Mynuddin, and M. Nabil, "AdaptEdge: Targeted universal adversarial attacks on time series data in smart grids," *IEEE Transactions on Smart Grid*, vol. 15, no. 5, p. 5072–5086, Sep. 2024.
- [8] M. Mynuddin, S. U. Khan, R. Ahmari, L. Landivar, M. N. Mahmoud, and A. Homaifar, "Trojan attack and defense for deep learning-based navigation systems of unmanned aerial vehicles," *IEEE Access*, vol. 12, p. 89887–89907, 2024.
- [9] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," 2019.
- [10] A. Niknam, H. K. Zare, H. Hosseinasab, and A. Mostafaeipour, "Developing an LSTM model to forecast the monthly water consumption according to the effects of the climatic factors in Yazd, Iran," *Journal of Engineering Research*, vol. 11, no. 1, p. 100028, Mar. 2023.
- [11] R. Qiu, Y. Wang, B. Rhoads, D. Wang, W. Qiu, Y. Tao, and J. Wu, "River water temperature forecasting using a deep learning method," *Journal of Hydrology*, vol. 595, p. 126016, Apr. 2021.
- [12] C. A. Bonilla, A. Zanfei, B. Brentan, I. Montalvo, and J. Izquierdo, "A digital twin of a water distribution system by using graph convolutional networks for pump speed-based state estimation," *Water*, vol. 14, no. 4, p. 514, Feb. 2022.
- [13] A. Menapace, A. Zanfei, M. Herrera, and B. Brentan, "Graph neural networks for sensor placement: A proof of concept towards a digital twin of water distribution systems," *Water*, vol. 16, no. 13, p. 1835, Jun. 2024.
- [14] Ambling Ingeniería y Servicios, S.L., "Ambling. Official Website," 2025, accessed: March 20, 2025. [Online]. Available: <https://www.ambling.es>
- [15] A. Rezvani, A. M. Saghiri, S. M. Vahidipour, M. Esnaashari, and M. R. Meybodi, *Recent Advances in Learning Automata*. Springer International Publishing, 2018.