



POLITECNICO DI TORINO

Master degree course in Digital Skills for Sustainable Societal
Transitions

Master Degree Thesis

Hydraulic Modeling and Machine Learning Solutions for Leak Monitoring in Municipal Water Distribution Networks

Relatori

Prof. Gianvito Urgese

Dr. Walter Gallego Gomez

Dr. Salvatore Tilocca

Candidato

Arman Moradi

Supervisore Aziendale

Elio Becchis (Fondazione DIG421)

December 2025

Abstract

Urban water distribution networks (WDNs) are essential infrastructures that ensure the continuous supply of clean and safe water to communities. Due to population growth, climate variability, and aging infrastructure, improving operational efficiency and minimizing losses have become critical objectives for water utilities. Among the major challenges faced by modern WDNs, leakages represent a significant source of inefficiency; in Italy, networks lose approximately 42% of the treated water before it reaches consumers. Such losses impose substantial economic costs, reduce system resilience, and hinder the sustainable management of water resources.

This thesis work involved the development of a hydraulic digital model and the evaluation of algorithms for leakage monitoring in two municipalities in the Province of Cuneo, Piedmont, Italy: Cavallermaggiore and Marene. The thesis was developed in collaboration with Alpi Acque, the water utility manager, and the support of Tesisquare and Fondazione DIG421. The study builds upon real infrastructure documentation, mainly GIS WDNs data, SCADA measurements, operational knowledge and records of historical leakages.

The thesis is divided into two complementary parts:

In Part I, we developed hydraulic digital models of Cavallermaggiore and Marene using the open source software EPANET. The models integrate reconstructed tank geometries, pump curves, well output data, district-level time patterns derived from SCADA flow measurements, yearly consumer-level demand, and rule-based controls that replicate real operational logic. The models were calibrated to match flow and pressure trends over a representative time period and were used to support the definition of districts in Cavallermaggiore and the development of Part II.

In Part II, we evaluated leakage monitoring methods based on machine learning techniques. First, we evaluated the use of a data-driven anomaly detection algorithm using real flow sensor data from the inlets of each district, to implement leakage detection. With this approach we are able to reach over 90% of recall of the historical leakages reported by the utility manager. Second, we used the validated hydraulic model to generate simulated leakage scenarios and pressure measurements, following the approaches described in the BattLeDIM competition. This data was then used to feed the LILA algorithm for leakage detection and localization.

By integrating a hydraulic digital model with machine learning techniques, this work provides a practical and scalable framework for monitoring leakage events in medium-sized municipal water distribution networks. The collaboration with Alpi Acque and the real-data-informed modelling of Cavallermaggiore and Marene demonstrate how hydraulic simulation, SCADA analytics, and machine learning solutions can support the transition toward smarter, more efficient, and more resilient water management.

Acknowledgments

I would like to express my sincere gratitude to my supervisor, Prof. Gianvito Urgese, for their invaluable guidance, constructive feedback, and continuous support throughout the development of this thesis. I am also grateful to Dr. Walter Gallego Gomez and Dr. Salvatore Tilocca for their insightful discussions and technical advice, which greatly contributed to shaping the methodology and improving the outcomes of this work.

Special thanks are extended to the technical staff and colleagues at Politecnico di Torino and Alpi Acque, the water utility manager, and the support of Tesisquare and Fondazione DIG421, who provided access to real network data and offered practical insights into the challenges of leakage management. Finally, I would like to thank my family and friends for their encouragement, patience, and unwavering support during this demanding but rewarding journey.

Contents

List of Figures	8
List of Tables	10
1 Introduction	11
2 Literature Review	15
2.1 Water Distribution Networks: Challenges and Water Losses	15
2.2 Hydraulic Modeling and Simulation	16
2.3 Model Calibration and Uncertainty Quantification	17
2.4 District Metered Areas and Network Sectorization	18
2.5 Leakage Detection Methods	18
2.5.1 Flow-Based Anomaly Detection	19
2.5.2 Pressure-Based Leakage Detection	19
2.6 Leakage Localization Approaches	20
2.6.1 Model-Based Localization	20
2.6.2 Data-Driven Localization	21
2.6.3 Mixed Approaches	21
2.7 The BattLeDIM Competition	21
2.8 The LILA Methodology	23
2.9 Sensor Placement and Network Design	25
2.10 Research Gaps	25
3 Study Area and Data Sources	29
3.1 Geographical and Operational Context	29
3.2 Infrastructure Data and GIS Master Records	30
3.3 Wells, Tanks, and Pumping Facilities	32
3.4 SCADA Measurements and Operational Data	34
3.5 Demand Data and Time Patterns	35
4 Part I: Hydraulic Model Development and Calibration	37
4.1 Static Stage: Steady State Model	38
4.1.1 Draft model from Technical Design	39

4.1.2	Draft model automatic corrections	39
4.1.3	Draft model manual corrections	39
4.1.4	Tank Modeling and Simulation	40
4.1.5	Well and Pump Modeling	40
4.1.6	Control Logic Implementation	41
4.2	Dynamic Stage: Extended-Period Simulation	42
4.2.1	Output Extraction and Analysis	44
4.3	Model Calibration	44
4.4	District Metered Area Definition for Cavallermaggiore	46
4.4.1	DMA Design Criteria	46
4.4.2	Simulation-Based Evaluation	47
5	Part II: Leakage monitoring algorithms	49
5.1	Flow based leakage detection	49
5.1.1	Overview and Motivation	49
5.1.2	DARTS Framework for Time-Series Forecasting and Anomaly Detection	50
5.1.3	Model Selection and Configuration	51
5.2	Pressure based leakage detection and localization	52
5.2.1	Simulation of leakages with BattleDIM approach for Cav- allermaggiore	52
5.2.2	LILA for detection	54
5.2.3	LILA for localization	55
6	Results from Leakage monitoring strategies	57
6.1	Leakage detection based on flow sensors	57
6.2	Leakage detection and localization based on pressure sensors	62
6.2.1	Initial evaluation on BattleDIM	62
6.2.2	Data generation	64
6.2.3	Leakage identification step	64
6.2.4	Leakage localization step	67
6.2.5	Results summary and discussion	69
7	Conclusion	73
	Bibliography	75

List of Figures

2.1	Schematic representation of the L-TOWN WDN	16
2.2	EPANET simulation example	17
2.3	Schematic of a Water Distribution Network	22
2.4	Simulated leak locations in a WDN	23
2.5	LILA methodology flowchart	24
3.1	Aerial view of Cavallermaggiore	29
3.2	Aerial view of Marene	30
3.3	GIS Master of the Cavallermaggiore WDN	31
3.4	GIS Master of the Marene WDN	32
3.5	Tank geometry parameters in EPANET	33
3.6	SCADA interface at Via Cuneo station	34
3.7	SCADA interface at Via Torino station	35
4.1	Hydraulic model construction and calibration workflow	38
4.2	Static and dynamic data to estimate nodes demand	43
4.3	Time Pattern example	43
4.4	An Example of pressure calibration	44
4.5	An Example of tank level calibration	45
4.6	Examples for tank level, flow and pressure sensor calibration	45
4.7	District Metered Area (DMA) definition	46
4.8	Histogram of pressure differentials before and after DMAs	47
4.9	Pressure calibration after (DMA) definition	48
6.1	Detected vs actual anomalies on BattLeDIM	58
6.2	ROC curve on the BattLeDIM dataset	59
6.3	Confusion matrix on the BattLeDIM 2018 dataset	59
6.4	Leak detection overview for Pellaverne, Marene	60
6.5	Nightly anomaly detection in Pellaverne	60
6.6	Leak detection overview for Marconi	61
6.7	Nightly anomalies in the Marconi district	61
6.8	Leak localization on BattLeDIM 2019 for MAS 276	62
6.9	All pressure sensors together	63

6.10	Statistics of simulated leakages	63
6.11	MAS for the Concentrico district, Cavallermaggiore	64
6.12	MAS for the Via Bra district, Cavallermaggiore	65
6.13	MAS for the Europa district, Cavallermaggiore	65
6.14	MAS for the Roma district, Cavallermaggiore	65
6.15	MAS for the Provinciale district, Cavallermaggiore	66
6.16	MAS for the Foresto district, Cavallermaggiore	66
6.17	MAS for the Madonna Del Pilone district, Cavallermaggiore	66
6.18	Correctly classified leak localization	67
6.19	Wrongly classified leak localization	68

List of Tables

6.1	Performance metrics of the leak detection model on the BattleDim dataset 2018.	58
6.2	Leak detection and localization results for Cavallermaggiore	70
6.3	Leak localization summary per district	71

Chapter 1

Introduction

Water distribution networks are essential for delivering clean and safe water to communities, yet they face mounting challenges from aging infrastructure, climate change, and growing demand. Among these challenges, water losses due to leakages represent one of the most pressing concerns for utilities worldwide. In Italy, the situation is particularly critical: approximately 42% of treated water is lost before reaching consumers [1], resulting in significant economic costs, wasted energy, and reduced system resilience. Addressing this challenge requires a combination of effective monitoring systems, predictive models, and automated detection algorithms.

The importance of leakage monitoring extends beyond economic considerations. Undetected leaks compromise water quality, increase operational costs through excess pumping, and undermine public confidence in water services. Traditional leakage detection methods, such as acoustic surveys and field inspections, are reactive and labor-intensive, making them unsuitable for continuous monitoring across large networks. As a result, there is growing interest in data-driven and model-based approaches that leverage real-time sensor data and hydraulic simulations to enable proactive leakage management.

Hydraulic modeling plays a central role in modern leakage management strategies. A well-calibrated hydraulic model enables utilities to simulate system behavior under various conditions, evaluate operational scenarios, and generate synthetic data for training detection algorithms. This becomes especially valuable when physical sensors are sparse or absent, as models can provide the hydraulic insights needed to interpret limited measurements and identify anomalies. Recent benchmark initiatives, such as the BattLeDIM competition [2], have demonstrated the effectiveness of combining hydraulic modeling with machine learning for leakage detection and localization in controlled synthetic environments.

However, a significant gap exists between research developments in synthetic benchmarks and practical applications in real operational networks. Most advanced leakage detection methods are evaluated using simulated data with idealized conditions, abundant sensor coverage, and perfect knowledge of network topology.

In contrast, real municipal water networks operate under substantial constraints: sensor deployments are limited by cost and maintenance requirements, hydraulic models are subject to calibration uncertainties, and leakage records are incomplete or delayed. While studies such as [3, 4, 5] have demonstrated promising results on real networks, comprehensive case studies that integrate flow-based detection, pressure-based localization, and real SCADA data on medium-sized municipal systems remain scarce.

This thesis addresses these gaps by developing and evaluating leakage monitoring methodologies on two real municipal water distribution networks: Cavallermaggiore and Marene, located in the Province of Cuneo (Piedmont, Italy). The work was conducted in collaboration with Alpi Acque S.p.A., the local water utility manager, with additional support from Tesisquare and Fondazione DIG421. The study builds on real operational assets, including GIS infrastructure data, SCADA measurements from flow and pressure sensors, historical leakage records documented by utility staff, and operational knowledge of system behavior. Both networks are characterized by complex topologies, multiple supply sources (wells and tanks), variable elevations, and heterogeneous consumption patterns typical of medium-sized Italian municipalities.

A key challenge addressed in this thesis is the constraint of limited sensor availability. Unlike synthetic benchmarks or large metropolitan networks where dense sensor coverage may be assumed, Cavallermaggiore and Marene have sparse instrumentation: flow meters are installed at district inlets, but pressure sensors are limited to a small number of critical nodes. This constraint reflects the operational reality faced by many small- and medium-sized utilities and requires adapted methodologies that maximize the value of limited measurements. The hydraulic model becomes essential in this context, providing virtual pressure observations and enabling scenario-based evaluation of detection and localization performance.

The remainder of this manuscript is organized as follows:

Chapter 2 offers an in-depth overview of the literature, covering hydraulic modeling and simulation, model calibration and uncertainty analysis, district metered area design, flow and pressure based leakage detection methods, the BattLeDIM benchmark, the LILA framework, and strategies for sensor placement.

Chapter 3 describes the study area and data sources, including the geographical and operational context of Cavallermaggiore and Marene, infrastructure data from GIS Master records, SCADA measurement systems, and historical leakage documentation provided by Alpi Acque.

Chapter 4 presents the hydraulic model development and calibration process, including network reconstruction from GIS data, demand allocation and pattern definition, tank and pump integration with control logic, extended-period simulation, calibration against SCADA flow and pressure measurements, validation, and the definition of district boundaries for Cavallermaggiore.

Chapter 5 evaluates leakage monitoring methods based on machine learning techniques, including flow-based anomaly detection applied to real SCADA measurements from district inlets, generation of synthetic leakage scenarios following the BattLeDIM methodology, and implementation of the LILA algorithm for pressure-based detection and localization under different sensor configurations.

Chapter 6 summarizes the main findings, discusses practical implications for water utilities managing networks under sparse sensor constraints, and suggests directions for future research in real-time deployment and SCADA integration.

By integrating real operational data, calibrated hydraulic models, and machine learning techniques, this thesis provides a practical and scalable framework for leakage monitoring in medium-sized municipal networks under realistic sensor constraints. The results demonstrate that even with limited instrumentation, effective leakage detection and localization can be achieved through the strategic of flow-based analytics and model-supported pressure analysis.

Chapter 2

Literature Review

2.1 Water Distribution Networks: Challenges and Water Losses

Water Distribution Networks (WDNs) are complex engineered systems composed of pipelines, storage tanks, pumps, valves, and nodes that deliver potable water from sources to consumers (Figure 2.1). Their performance depends on the ability to maintain adequate pressure levels, ensure water quality, and minimize hydraulic and operational inefficiencies. Over the past decades, WDNs worldwide have experienced increasing pressures due to demographic growth, climate variability, and aging infrastructure [6, 7].

Real water losses remain one of the most critical challenges for water utilities. These losses primarily originate from leakages caused by pipe deterioration, faulty joints, corrosion, ground movement, and pressure fluctuations [9, 10]. The International Water Association (IWA) has developed standardized terminology and performance metrics to quantify and manage losses across different contexts [7, 11]. In many regions, non-revenue water exceeds 30% of total production, imposing substantial economic costs and reducing system resilience [9].

In Italy, the situation is particularly severe, with annual leakage rates exceeding 42% in many municipalities, reinforcing the urgent need for advanced monitoring and modeling strategies [1]. Excessive losses increase utility operating costs, reduce network resilience, complicate long-term planning for sustainable water management, and contribute to resource depletion. International studies highlight that pressure control, sectorization through District Metered Areas (DMAs), and continuous monitoring play essential roles in reducing leakage levels [12, 13, 14].

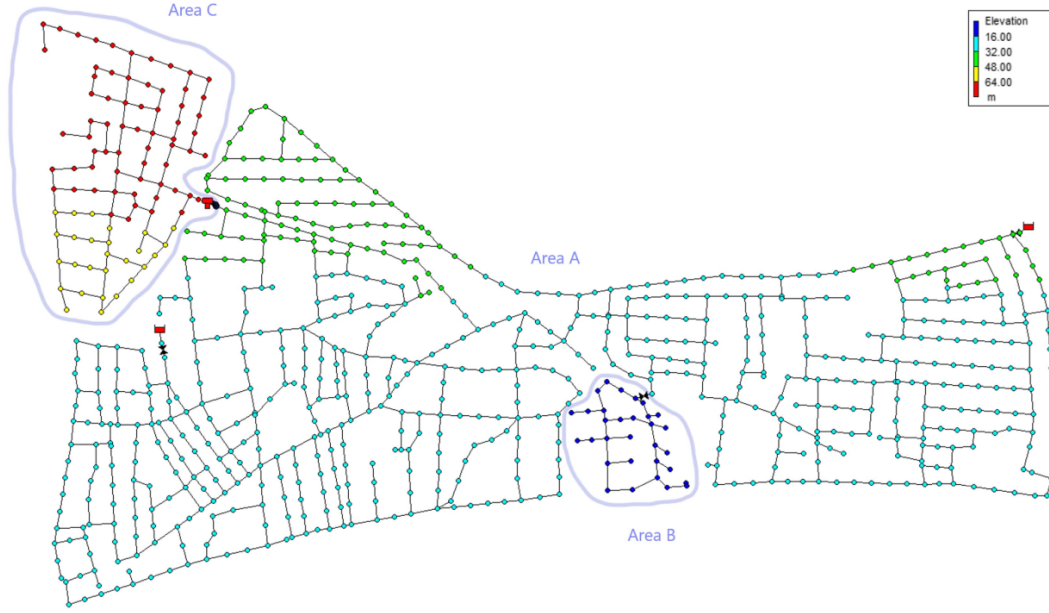


Figure 2.1: Schematic representation of the L-TOWN Water Distribution Network showing main components: pipelines, storage tanks, pumps, valves, and demand nodes[2, 8].

2.2 Hydraulic Modeling and Simulation

Hydraulic modeling is widely adopted by water utilities for design, planning, and real-time decision support. Software tools such as EPANET, WaterGEMS and InfoWorks WS enable simulation of pressure, flow, and tank dynamics across extended periods, allowing analysis of operational scenarios and fault conditions [15, 16]. A calibrated hydraulic model is fundamental for reproducing system behavior, supporting design decisions, and enabling leakage management strategies [17].

EPANET (Figure 2.2), developed by the U.S. Environmental Protection Agency, is the most widely used open-source hydraulic modeling platform for water distribution systems. It implements the demand-driven analysis (DDA) approach by default, where nodal demands are assumed to be satisfied regardless of pressure availability. However, under abnormal operating conditions such as leakages or pipe failures, pressure-driven analysis (PDA) becomes more appropriate, as it accounts for the relationship between available pressure and actual delivered demand [18, 19].

Modeling workflows generally include reconstruction of network topology from GIS data, assignment of nodal demands and time patterns, characterization of pumps, valves, and tanks, and validation of model predictions against SCADA measurements [20]. Extended-period simulation (EPS) allows for the analysis of

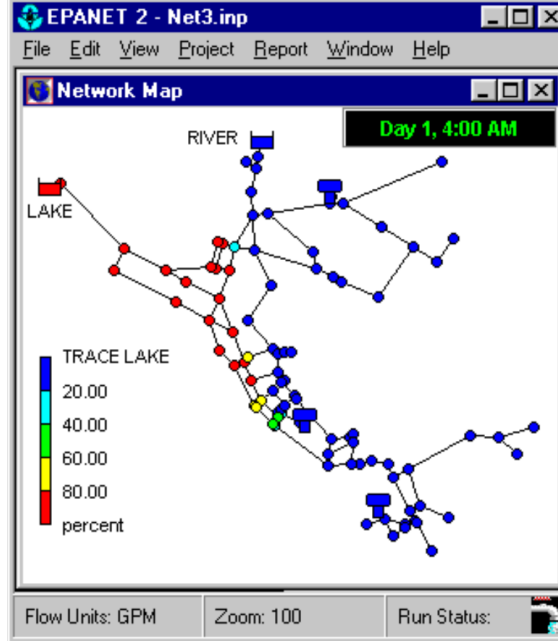


Figure 2.2: EPANET simulation example. (Image from EPANET’s Wikipedia).

time-varying behavior, including daily demand cycles, tank level fluctuations, and pump scheduling [16]. This capability is particularly important when evaluating the hydraulic signatures of leakages, which may evolve gradually over hours or days.

2.3 Model Calibration and Uncertainty Quantification

Model calibration is the process of adjusting uncertain model parameters to minimize discrepancies between simulated and observed hydraulic data. Key parameters subject to calibration include pipe roughness coefficients, nodal demands, pump characteristics, and control logic [20]. Calibration techniques range from manual trial-and-error adjustments to advanced optimization approaches based on evolutionary algorithms, surrogate modeling, or Bayesian inference [21, 22, 23].

Recent works emphasize the role of uncertainty quantification in improving model robustness. Sources of uncertainty include measurement errors in SCADA data, inaccuracies in GIS records, variability in consumer demand, and simplifications in the hydraulic representation [24, 25]. Properly accounting for these uncertainties is crucial when using models to support critical decisions such as leak localization or operational optimization.

Demand estimation is a particularly challenging aspect of model calibration. Real-time demand estimation methods have been developed to dynamically adjust

nodal demands based on observed flow and pressure measurements [26]. These methods are valuable for improving model accuracy under varying operational conditions and for detecting anomalies that may indicate leakages or unauthorized consumption.

2.4 District Metered Areas and Network Sectorization

District Metered Areas (DMAs) are hydraulically isolated zones within a distribution network, typically defined by closing boundary valves and installing flow meters at the inlets (Figure 2.1). DMAs enable water utilities to monitor consumption at a finer spatial resolution, identify abnormal flow patterns, and prioritize areas for leak detection campaigns [27, 14].

The creation of DMAs supports leakage management in several ways. First, it reduces the search space for leak localization by narrowing the focus to a specific district. Second, it facilitates minimum night flow (MNF) analysis, where flow measurements during low-demand periods (typically 2:00 AM to 4:00 AM) are used to estimate background leakage levels [13]. Third, it allows for controlled pressure management within each zone, reducing stress on aging pipes and slowing the development of new leaks [28].

However, DMA design requires careful consideration of hydraulic constraints, including pressure requirements, redundancy, and emergency supply routes. Recent studies have explored automated algorithms for optimal DMA boundary definition, balancing monitoring benefits with operational flexibility [29]. In this thesis, the hydraulic model of Cavallermaggiore was used to support the definition of district boundaries and to evaluate the impact of sectorization on network performance.

2.5 Leakage Detection Methods

Leakage detection refers to the process of identifying the presence and timing of anomalous water losses within a distribution network, distinguishing leakage events from normal operational variability and measurement noise [5, 9]. Leakage detection is traditionally performed through acoustic surveys, step testing, ground-penetrating radar, and field inspections. While effective for locating leaks once they are suspected, these methods are labor-intensive, expensive, and do not scale to large network areas [9, 10]. SCADA-based monitoring provides real-time flow and pressure data but is often limited by sparse sensor availability and measurement noise. As a result, utilities increasingly rely on data-driven and model-based methods to enhance situational awareness and enable proactive leak management [30].

Leakage events typically manifest as deviations in flow or pressure compared to expected behavior. Detection approaches can be broadly categorized into flow-based methods and pressure-based methods, each with distinct advantages and limitations depending on the available instrumentation and network characteristics.

2.5.1 Flow-Based Anomaly Detection

Flow-based methods exploit trends at district or subdistrict inlets, where sensors are more commonly installed. These methods analyze time series of flow measurements to detect abnormal consumption patterns or sudden increases in minimum night flow, which may indicate the presence of leakages [31, 3].

Anomaly detection models based on statistical control charts, forecasting residual analysis, and machine learning techniques have been successfully applied to SCADA flow data. Statistical process control (SPC) methods, such as CUSUM (Cumulative Sum) charts, detect shifts in the mean or variance of a monitored signal [32, 33]. Time-series forecasting approaches, including ARIMA, exponential smoothing, and modern machine learning frameworks such as Prophet [34] and DARTS [35], generate predictions of expected flow under normal conditions. Persistent deviations between observed and predicted values trigger alarms that may indicate leakage events.

Machine learning techniques, including artificial neural networks (ANN), support vector machines (SVM), and long short-term memory (LSTM) networks, have also been explored for flow anomaly detection [36, 37]. These methods can capture complex nonlinear relationships and seasonal patterns in consumption data, improving detection performance in networks with heterogeneous demand profiles.

Recent literature shows that flow-based methods can successfully identify historical leakages when supported by reliable SCADA measurements at DMA inlets [4]. This aligns with the approach adopted in Part II of this thesis, where a data-driven anomaly detection algorithm is applied to real SCADA flow measurements to detect historical leakage events in the Cavallermaggiore network.

2.5.2 Pressure-Based Leakage Detection

Pressure-based leakage detection relies on the principle that leaks cause localized pressure drops throughout the network. The magnitude and spatial pattern of these pressure changes depend on leak size, location, network topology, and operating conditions [38, 39]. However, pressure sensors are typically sparse in operational WDNs due to installation costs, maintenance requirements, and data transmission constraints.

To overcome these limitations, recent research has explored the combination of hydraulic models and machine learning to amplify the value of limited pressure measurements. Pressure residuals, defined as the difference between observed and

expected (model-predicted) pressures, serve as indicators of abnormal conditions [5]. Statistical change detection algorithms, such as CUSUM, can be applied to these residuals to identify the onset of leakage events.

Transient-based methods, which analyze pressure wave propagation following sudden changes in flow or valve operations, have also been investigated for rapid leak detection. These methods offer the potential for near-instantaneous detection but require high-frequency pressure measurements and specialized instrumentation [10]. In practice, most operational systems rely on quasi-steady-state pressure measurements collected at intervals ranging from minutes to hours.

2.6 Leakage Localization Approaches

While leakage detection identifies when and whether a leak is present, leakage localization aims to determine where the leak is occurring within the network, requiring the spatial inference of leak positions from sparse sensor measurements [40, 41]

Localizing leaks is significantly more challenging than detecting them. Pressure gradients propagate across the network in complex nonlinear patterns, making inversion of leakage-induced signatures an ill-posed problem. Small leaks may produce weak and spatially diffuse pressure signatures, particularly in highly interconnected networks with significant background uncertainty [41].

The main approaches are commonly described in the literature: model-based localization relying on hydraulic simulations and residual matching, data-driven localization exploiting sensor correlations and machine learning classifiers, and hybrid methods that combine both paradigms [40, 5].

2.6.1 Model-Based Localization

Model-based localization involves comparing observed pressure variations with simulated leakage patterns at candidate nodes. For each node in the network (or within a reduced search space such as a DMA), a synthetic leak is introduced in the hydraulic model and the resulting pressure distribution is computed. The simulated pressure residuals are then compared against observed residuals using correlation metrics, Euclidean distance, or other similarity measures [42].

Bayesian classifiers have been explored to formalize the probabilistic inference process, combining prior knowledge (e.g., pipe age, historical leak frequency) with likelihood functions derived from hydraulic simulations [43]. However, these methods require significant computational effort when applied to large networks with hundreds or thousands of candidate nodes.

Search-space reduction strategies, such as coarse-to-fine localization or Most Affected Sensor (MAS)-based filtering, are essential for improving computational

efficiency and localization [41]. In this thesis, the MAS approach is used to inform the leakage localization step, restricting the search space to nodes identified as hydraulically sensitive to the detected leakage event.

2.6.2 Data-Driven Localization

Data-driven methods use clustering, dimensionality reduction, or supervised machine learning to infer leak locations from sensor patterns without explicit hydraulic modeling [44, 40]. These methods typically require large training datasets of labeled leakage events, which are often generated synthetically using hydraulic models.

Support vector machines (SVM), artificial neural networks (ANN), and ensemble classifiers such as random forests have been applied to leak localization problems. These methods can capture complex nonlinear relationships between sensor measurements and leak locations, potentially improving performance in networks with heterogeneous topology or time-varying demand patterns.

However, data-driven methods are sensitive to the quality and representativeness of the training data. Overfitting, poor generalization to unseen scenarios, and lack of interpretability remain challenges.

2.6.3 Mixed Approaches

Recognizing the complementary strengths and limitations of purely model-based and purely data-driven methods, hybrid methodologies have emerged that integrate both paradigms to leverage the physical interpretability of hydraulic models alongside the pattern recognition capabilities of machine learning techniques [5, 38, 40]. Common strategies include sequential coupling, where data-driven methods reduce the search space before model-based refinement [40], and parallel fusion through ensemble voting or Bayesian methods [38, 41]. A particularly important approach uses hydraulic simulators like EPANET to generate synthetic training datasets for machine learning classifiers, transferring knowledge to operational networks where real labeled data is scarce [8, 45]. While promising, challenges remain in tuning, computational overhead, and interpretability [41].

2.7 The BattLeDIM Competition

The Battle of the Leakage Detection and Isolation Methods (BattLeDIM) was a benchmark competition organized in 2020 to evaluate and compare leakage detection and localization algorithms under controlled conditions [2, 8]. The competition was based on a water distribution network called L-TOWN, significantly modified and extended to create a synthetic benchmark suitable for algorithm testing, (Figure 2.3).

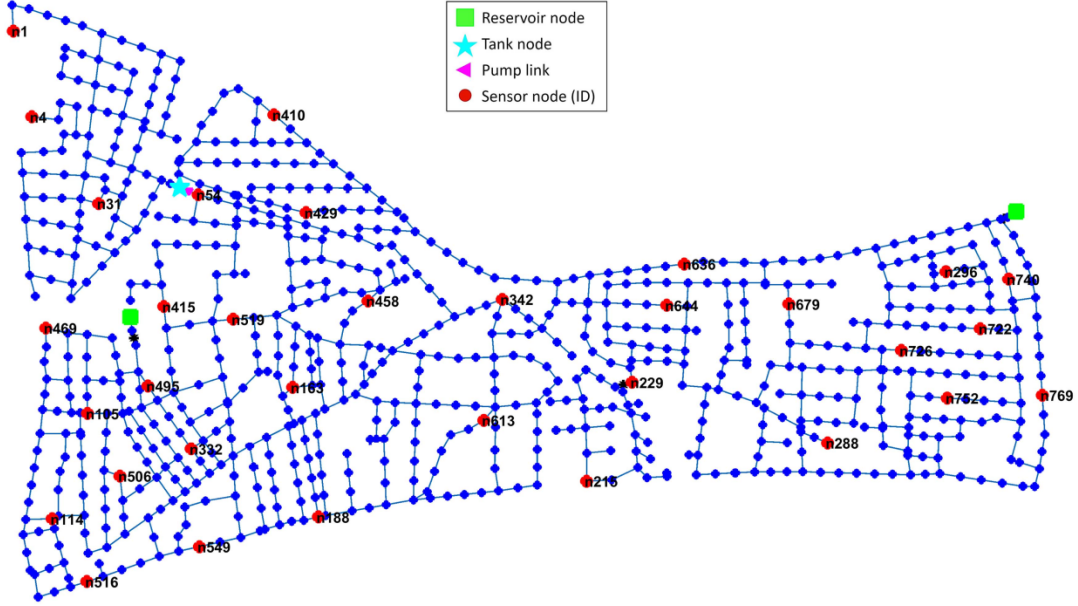


Figure 2.3: Schematic representation of a Water Distribution Network showing main components: reservoir, tank, pumps, sensor nodes[2, 8].

The BattLeDIM team created an EPANET model, intended to represent the WDN. Multiple leakage events with varying sizes, locations, and temporal characteristics (incipient and abrupt leaks) were simulated starting from this model. Figure 2.4).

The following material is given to the competitors:

- **Network Model:** A “noisy” EPANET model representing an imperfect approximation of the true network, generated by adding uncertainty to pipe roughness coefficients and nodal demands of the “clean” model. This reflects the real-world condition that network models are never perfectly calibrated. The true network model was never revealed to participants.
- **SCADA Measurements:** Synthetic time-series data including pressure measurements at selected nodes, flow measurements at pipes and DMA inlets, tank levels, and AMR (Automated Meter Reading) consumption data at selected nodes, obtained by running the true network with the simulated leakage scenarios.
- **Historical data:** of simulated leakages present in the network, so the competitors can train and evaluate their models. This set of leakages is different from the one the competitors need to identify.

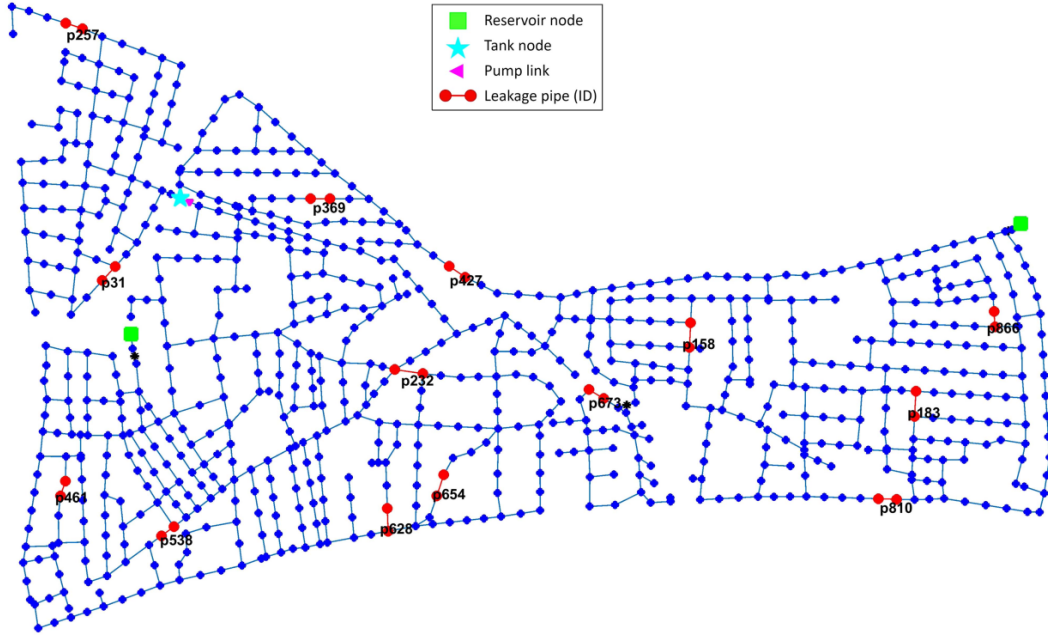


Figure 2.4: Schematic representation of a Water Distribution Network showing simulated leakage locations.[2, 8].

The objective of the competitors is to develop algorithms capable of detecting leakage events from SCADA time-series data and accurately localizing the leak positions within the network using only the imperfect hydraulic model and sparse sensor measurements provided [2, 8].

The BattLeDIM dataset has become a widely used benchmark for developing and validating leakage detection and localization algorithms. It demonstrated the feasibility of combining hydraulic simulation with machine learning to create synthetic training data for scenarios that may be rare or difficult to observe in operational networks [8].

Participants were evaluated using a utility-oriented performance framework that rewards timely and accurate leak detection and spatial localization within predefined tolerances, while penalizing false alarms and delayed responses, in order to reflect the operational impact on real water utilities [2, 8].

2.8 The LILA Methodology

Among the top-performing methods in BattLeDIM, the LILA (Leakage Identification and Localization Algorithm) methodology proposed by Daniel et al. [45]

introduced a sequential detection and localization framework based on linear regression models and pressure residual analysis. LILA achieved 3rd (out of 18) place (Third Place Award) in the BattLeDIM 2020 competition. However, the original implementation relies on manual selection of sensor pairs for detection and expert definition of MAS for localization. The LILA methodology consists of two main steps: (Figure 2.5):

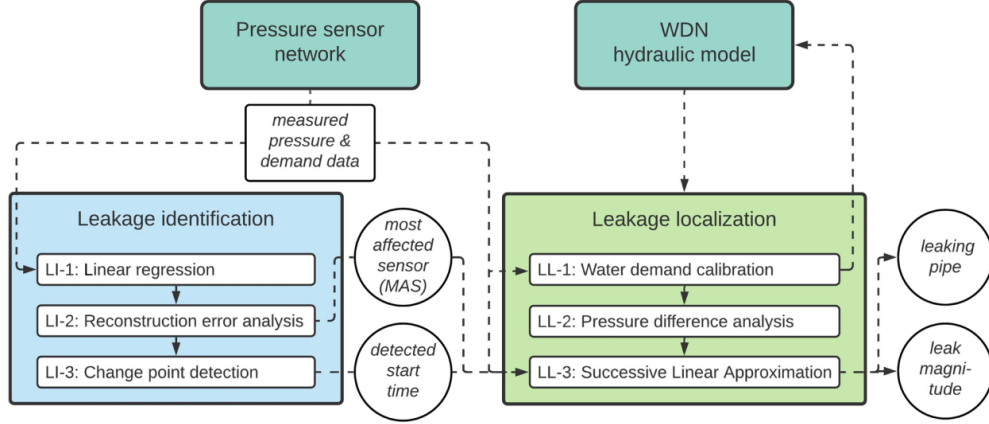


Figure 2.5: Flowchart of LILA, two-step method for leakage identification and localization.

Stage 1: Leakage Detection LILA detects leakages through pairwise linear regression analysis of pressure sensor data. The detection stage includes:

- **Regression Modeling:** LILA constructs pairwise linear regression models between pressure sensors under normal (leak-free) conditions. For each sensor pair (i, j) , a regression model $P_j = \alpha_{ij} + \beta_{ij}P_i$ is trained using historical data.
- **Residual Calculation:** During operational monitoring, the residual error $e_j = P_j^{\text{obs}} - P_j^{\text{pred}}$ is computed for each sensor. The Median Residual Error (MRE) aggregates residuals across sensors to provide a robust leakage indicator.
- **CUSUM Detection:** A CUSUM chart is applied to the MRE signal to detect statistically significant deviations from normal behavior. When the CUSUM statistic exceeds a predefined threshold, a leakage event is flagged.

Stage 2: Leakage Localization Once a leak is detected, LILA performs model-based localization:

- **Most Affected Sensor (MAS):** The sensor with the largest residual during the detected event is identified as the Most Affected Sensor. This sensor is assumed to be hydraulically closest to the leak.

- **Candidate Node Selection:** A subset of candidate leak locations is defined based on hydraulic proximity to the Most Affected Sensor, reducing the search space from all network nodes to a manageable subset.
- **Sensitivity Analysis:** For each candidate node, a synthetic leak is simulated using the available hydraulic model. The simulated pressure residuals at all sensor locations are compared against observed residuals to rank candidate locations.
- **Correlation Matching:** The candidate leak location with the highest correlation between simulated and observed residual patterns is selected as the most likely leak location.

2.9 Sensor Placement and Network Design

Optimal sensor placement is a critical component of effective leakage detection and localization. The goal is to determine the minimum number of sensors and their locations such that detection sensitivity and localization accuracy are maximized, subject to cost and operational constraints [46, 47].

Sensor placement strategies can be broadly classified into model-based and data-driven approaches. Model-based methods rely on hydraulic simulation to evaluate the sensitivity of each potential sensor location to leakages at different nodes [48]. Sensitivity matrices, which quantify the pressure change at each sensor due to a unit leak at each node, are often used to formulate optimization problems that maximize coverage or minimize uncertainty.

Data-driven methods use historical SCADA data and machine learning to identify sensor configurations that maximize information gain or detection performance [49]. Evolutionary algorithms, such as genetic algorithms (GA) and multi-objective optimization frameworks, are commonly employed to search the large combinatorial space of possible sensor configurations [48]. In practice, sensor placement must balance competing objectives, including detection sensitivity, localization accuracy, redundancy, and cost. Recent studies emphasize the importance of considering model uncertainty and demand variability when designing sensor networks [47]. In this thesis, a fixed sensor configuration of 10 pressure sensors is used, reflecting a realistic constraint for medium-sized municipal networks.

2.10 Research Gaps

The literature highlights the increasingly important role of calibrated hydraulic models, SCADA analytics, and machine learning techniques in modern leakage management. Flow-based anomaly detection, synthetic hydraulic simulation, and

regression-based pressure analysis each contribute to improving monitoring capabilities in networks with limited instrumentation [4, 45].

However, several research gaps remain:

- **Real-world validation:** While synthetic benchmarks such as BattLeDIM provide controlled testing environments, validation on real operational networks with actual SCADA data and documented leakage events remains limited. Notable exceptions include applications on real networks such as those reported in [3, 4, 5], yet comprehensive case studies combining detection and localization on medium-sized municipal systems are still relatively scarce.
- **Limited data availability:** Many water distribution networks, particularly in small to medium-sized municipalities, lack comprehensive historical SCADA data, accurate GIS records, or documented operational parameters necessary for traditional model calibration. Developing reliable hydraulic models under these data-scarce conditions requires alternative calibration strategies that can work with minimal measurements [20, 21]. The challenge of constructing and validating hydraulic models when baseline operational data is limited or entirely absent represents a significant practical barrier to implementing model-based leakage detection and localization methodologies in real-world settings.
- **Sparse sensor coverage challenges:** Many existing methods assume relatively dense sensor deployment or focus on networks where instrumentation is abundant. However, real operational constraints often result in sparse sensor coverage, where only a limited number of pressure sensors are available across the network. Evaluating detection and localization performance under these constrained conditions, particularly in networks where sensor-to-node ratios are low, represents a practical challenge that requires further investigation.

The methodologies adopted in this thesis address these gaps by:

1. Performing an analysis of the limited available data in a real-world WDN, and how this limited data can be exploited to build reliable software components.
2. Developing a calibrated hydraulic model of real-world WDS based on real GIS and SCADA data.
3. Applying flow-based anomaly detection to historical leakage records from district inlet flow meters.
4. Implementing a simplified pressure-based localization approach inspired by LILA principles, adapted to the constraints of limited source code availability.
5. Evaluating detection and localization performance under realistic noise levels and sensor configurations.

These contributions provide a practical and scalable framework for leakage management in medium-sized municipal water distribution networks, building on established methodologies such as BattLeDIM and LILA while adapting them to the constraints and opportunities of real-world operational environments.

Chapter 3

Study Area and Data Sources

3.1 Geographical and Operational Context

The study area comprises the municipal water distribution networks of Cavallermaggiore and Marene (Figures 3.1 and 3.2), located in the Province of Cuneo (Piedmont, Italy). Both municipalities are managed by Alpi Acque S.p.A. and form part of a broader regional supply system characterized by a combination of groundwater abstraction, local storage tanks, and pressure management through boosting stations. The two networks display heterogeneous elevation profiles, ranging from low-lying residential sectors to areas located at higher altitudes, which influences pressure levels and pumping requirements.

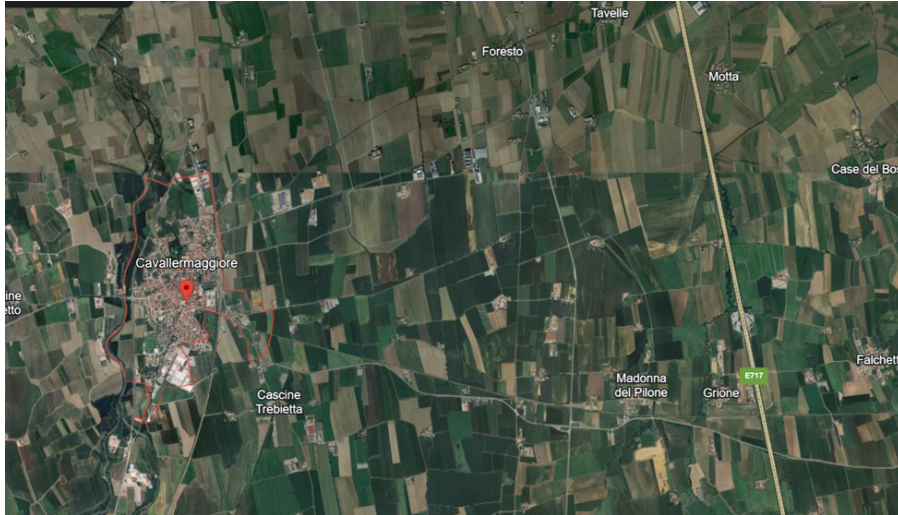


Figure 3.1: Aerial view of the Cavallermaggiore municipality showing the urban layout and water distribution infrastructure coverage area.



Figure 3.2: Aerial view of the Marene municipality showing the urban layout and water distribution infrastructure coverage area.

Cavallermaggiore is supplied mainly by groundwater wells and a set of storage tanks that regulate daily fluctuations in demand. The distribution network serves residential, commercial, and agricultural consumers and includes multiple pressure zones defined by local topography and operational constraints. Marene exhibits a similar configuration but relies on a distinct set of wells and tanks. The hydraulic behavior of both systems is governed by daily consumption cycles, tank level variations, and operational pump controls implemented by the utility.

3.2 Infrastructure Data and GIS Master Records

The hydraulic models developed in this thesis are based on the official GIS Master datasets provided by Alpi Acque through their technical design platform [50]. These datasets include detailed information on pipelines, nodes, valves, tanks, pumps, and wells, together with metadata describing installation dates, material types, pipe diameters, and nominal pressures, (Figures 3.3 and 3.4). The GIS layers were exported in vector format and processed to reconstruct the network topology required for hydraulic simulation.

Key attributes extracted from the GIS Master records include:

- pipe connectivity, length, diameter, roughness classification, and material;
- node elevations derived from survey data integrated in the GIS platform;



Figure 3.3: GIS Master representation of the Cavallermaggiore water distribution network showing pipes, nodes, tanks, and pumps extracted from the utility database.

- metadata for storage tanks (geometry, overflow elevation, minimum and maximum operational levels);
- pump characteristics and station layout;
- district boundaries and flow meter locations.

These datasets provided the structural baseline for generating the EPANET input files of Cavallermaggiore and Marene. Missing or inconsistent entries were addressed through cross-checks with operational maps, engineering drawings, and discussions with the utility’s technical staff.

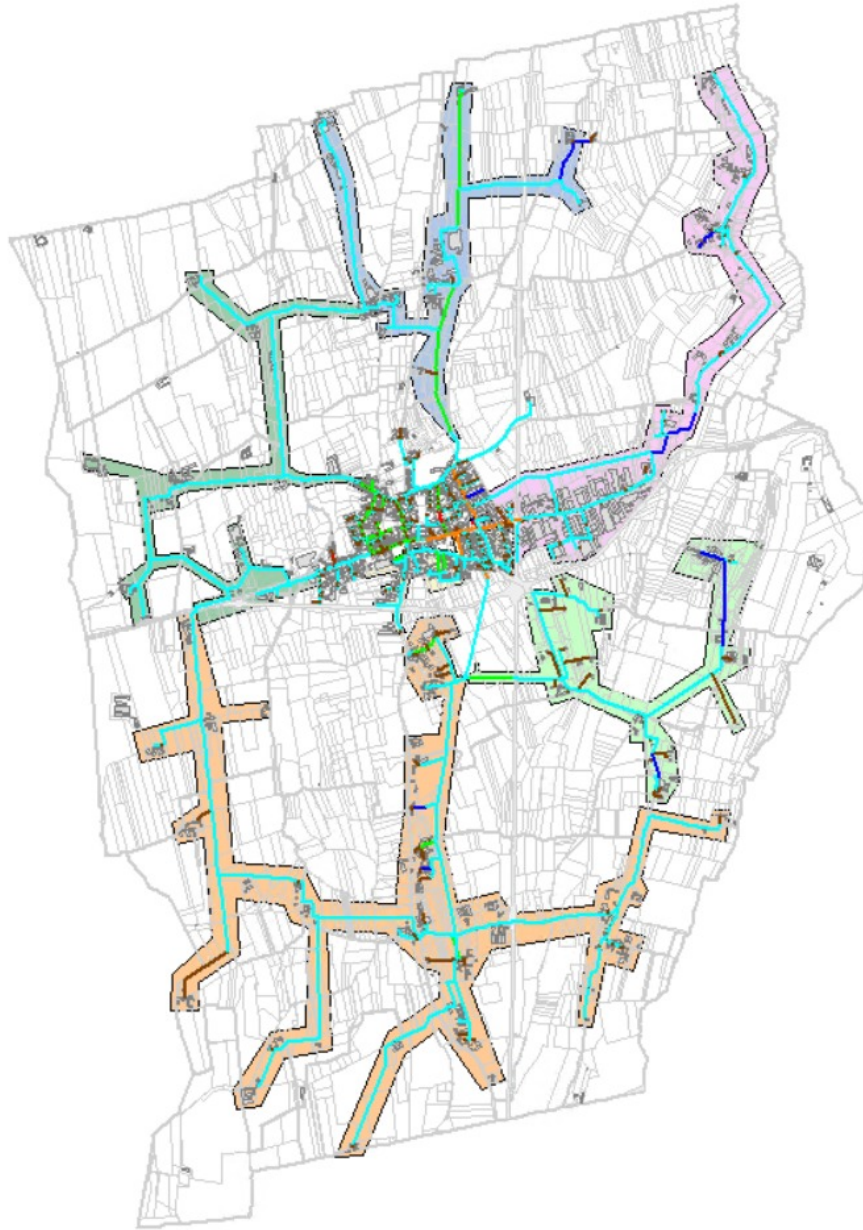


Figure 3.4: GIS Master representation of the Marene water distribution network showing pipes, nodes, tanks, and pumps extracted from the utility database.

3.3 Wells, Tanks, and Pumping Facilities

Groundwater wells constitute the primary source of supply for both systems. Operational data from the utility indicate average well output, daily variability, and pump activation cycles. Tanks play a central role in balancing intra-day demand

and stabilizing pressure across sectors. Their geometries were reconstructed using capacity curves, survey measurements, and information provided by the utility company (Figure 3.5).

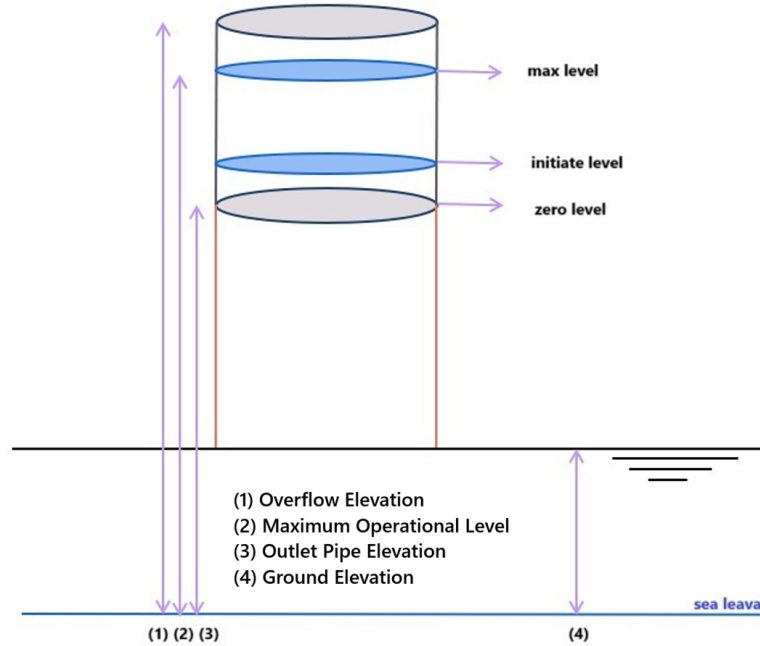


Figure 3.5: Schematic representation of tank geometry parameters used in EPANET modeling, including overflow elevation, maximum operational level, outlet pipe elevation, and ground altitude.

Pump curves, tank level constraints, and control logic were integrated into the hydraulic model to reproduce real system behavior. This included:

- pump characteristic curves derived from manufacturer data or field measurements;
- rule-based controls based on tank levels or time schedules;
- interconnections between supply zones used during specific operating conditions.

These parameters are essential for accurately simulating the hydraulic response of the networks over an extended period.

3.4 SCADA Measurements and Operational Data

SCADA measurements were provided by Alpi Acque through the IDEA platform (Figures 3.6 and 3.7). These data were used for demand pattern extraction, calibration, and performance evaluation. The available data consist of:

- flow at district inlets and main supply branches;
- pressures at some of the tank outlets, pumping stations, and selected network nodes;
- flow at well pumps output and pump status logs;
- tank levels;

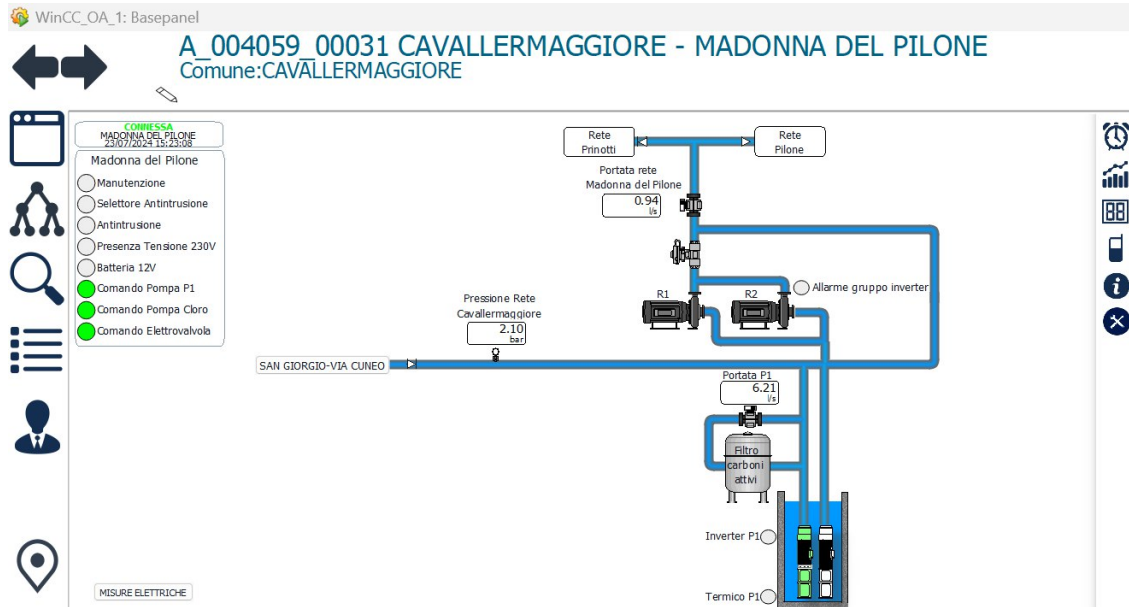


Figure 3.6: SCADA visualization interface for Cavallermaggiore Via Cuneo monitoring station showing real-time pressure, flow, and tank level measurements used for model calibration and operational monitoring.

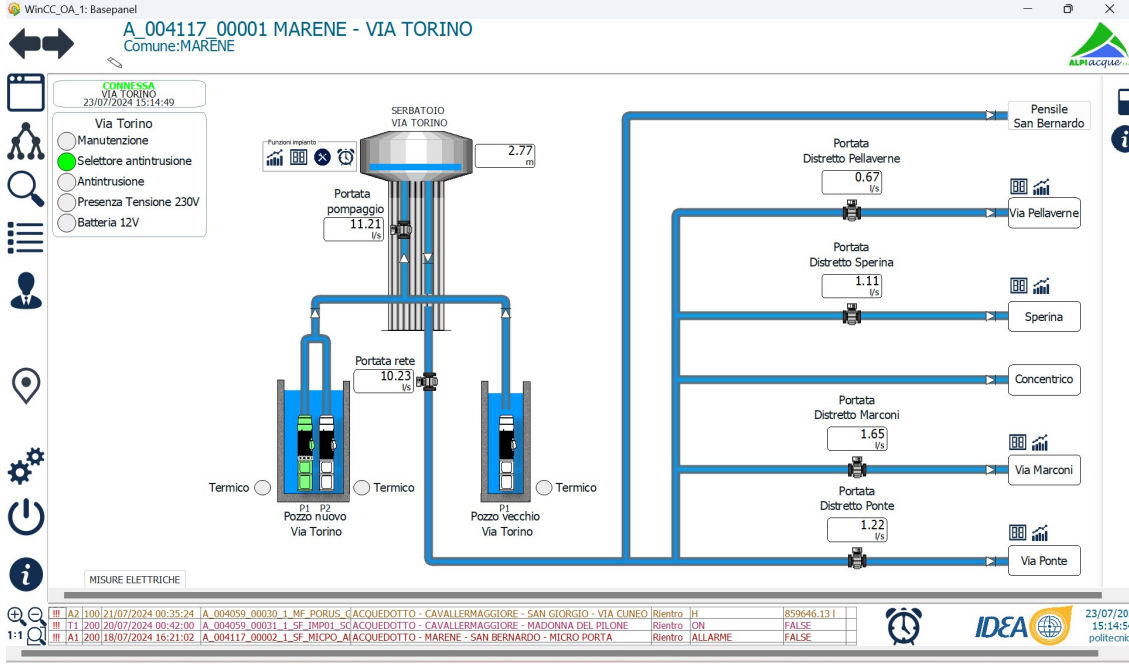


Figure 3.7: SCADA visualization interface for Marene Via Torino monitoring station showing real-time pressure, flow, and tank level measurements used for model design, calibration and operational monitoring.

3.5 Demand Data and Time Patterns

Yearly consumption data were obtained from the utility’s billing system, which records metered consumption at individual household and commercial connection points (contatori). These annual volumes were spatially distributed across the network model at the node level based on service connections documented in the GIS Master platform, ensuring that each demand node reflected the consumption of its associated users.

The temporal variability of consumption was captured through time patterns derived from SCADA flow sensors at the district inlets and reconciled with the yearly consumption records. This integration of SCADA-derived patterns with yearly consumer data ensured consistency between the spatial distribution of demands and the temporal variability observed in real operation, forming the basis for the extended-period simulations and calibration described in the next chapter.

Chapter 4

Part I: Hydraulic Model Development and Calibration

This chapter describes the development, simulation, and calibration of the hydraulic models for Cavallermaggiore and Marene. The models were built using EPANET, integrating GIS Master infrastructure data, SCADA measurements, and operational knowledge provided by Alpi Acque. The workflow included network topology reconstruction, demand allocation, tank and pump integration, control logic implementation, extended-period simulation, and calibration against observed tank level and pressure measurements. The workflow, summarized in Figure 4.1, is divided into three main stages:

Static stage in which the static input data, that is, data that does not change often, such as network topology and plant information, is used to build a **Steady** state model. This stage requires both automatic and manual steps, and requires an expert intervention. The stage is performed only when there are changes to the static data.

Dynamic stage in which the dynamic input data, that is, data that changes every day, such as SCADA flow measurements, is used to augment the Steady state model, resulting in a **Extended** period simulation model. This stage is automatic and it is performed at the beginning of each day.

Calibration stage in which the Steady state model is calibrated by using one Extended period simulation model and the calibration input data, composed of dynamic data that was not used during the dynamic stage, such as tank levels and pressure measurements. This stage requires expert intervention, and is performed every time the steady model changes, or periodically to ensure the correctness of the model.

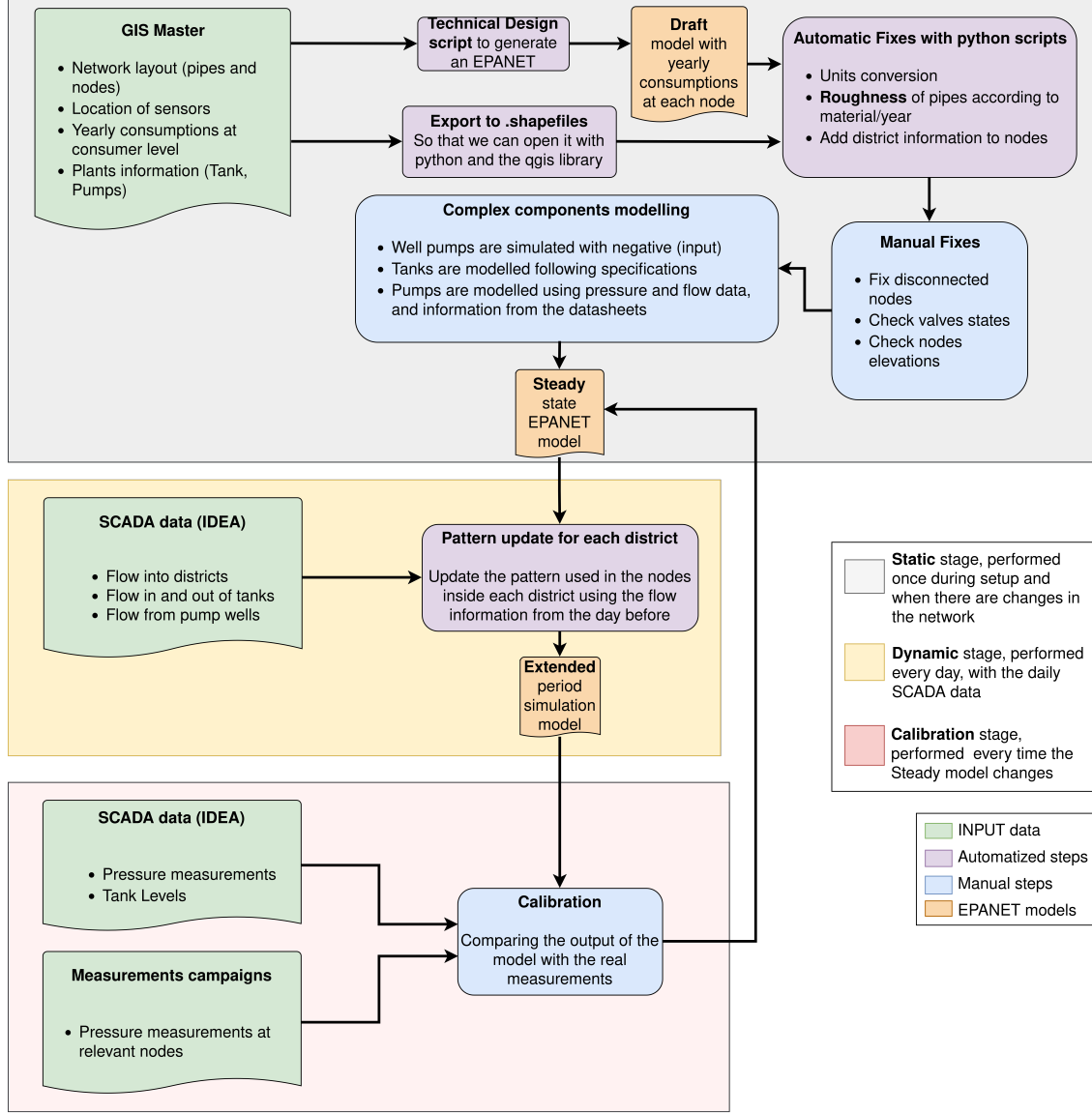


Figure 4.1: Workflow used for hydraulic model constructions and calibration

4.1 Static Stage: Steady State Model

The foundation of both hydraulic models is the network topology extracted from GIS Master records provided by Alpi Acque through their technical design platform [50]. The GIS datasets contain comprehensive information about network components, including pipes, nodes, valves, tanks, pumps, and wells considering EPANET conceptions [6, 16].

4.1.1 Draft model from Technical Design

Technical design provided a draft EPANET model automatically obtained from the GIS data. This draft model included:

- Start and end node coordinates (used to define network connectivity), with elevation information.
- Pipe length, nominal diameter, material type and installation date
- Yearly consumption volumes from the billing database—recorded at individual household and commercial connection points (*contatori*)—were spatially assigned to network nodes using service-connection information available in the GIS Master platform. This procedure establishes the base (static) demand distribution across the network.

Although this draft EPANET model contains valuable information, it is not complete and it is not even possible to perform a simulation from it. However it served as a starting point used to build the working models.

4.1.2 Draft model automatic corrections

The following corrections were necessary, and were done automatically with a Python script:

- Hydraulic units were converted to be compatible with SCADA units.
- Pipe roughness coefficients were assigned based on material type and age using standard reference values from the literature [16].
- District information was added to each node, in the *Pattern* field. This is used later to add the dynamic information.

4.1.3 Draft model manual corrections

- Elevation consistency was verified by comparing the EPANET values with online topography information sources for some selected nodes.
- Nodes disconnected in the EPANET model were properly connected.
- Valves documented in the GIS dataset were integrated into the model according to their functional role. District boundary valves were represented with closed status to reflect the sectorized configuration.

4.1.4 Tank Modeling and Simulation

Storage tanks play a central role in the hydraulic behavior of both networks, regulating pressure fluctuations, balancing supply and demand, and providing operational flexibility. Accurate representation of tank geometry, operating levels, and inlet/outlet configurations is essential for reproducing real system dynamics.

Tank Geometry Reconstruction

Tank geometries were reconstructed using survey measurements, and operational data provided by the utility. Tanks in Cavallermaggiore and Marene are cylindrical, allowing straightforward calculation of diameter from known volume and height:

$$D = 2\sqrt{\frac{V}{\pi \cdot H}} \quad (4.1)$$

where D is the tank diameter, V is the total volume, and H is the height. As shown in Figure 3.5, key elevation parameters include ground altitude, outlet pipe elevation, maximum operational level, and overflow elevation.

These parameters were carefully extracted from technical documentation and cross-checked with SCADA tank level measurements.

Inlet and Outlet Pipe Configurations

Tank inlet pipes were configured to connect supply sources (wells or upstream zones) to the tank. Outlet pipes distribute water from the tank to the downstream network. In cases where tanks are equipped with level-dependent controls or multiple inlet/outlet pipes serving different pressure zones, these configurations were explicitly modeled using EPANET's pipe and valve elements.

4.1.5 Well and Pump Modeling

Groundwater wells constitute the primary water source for both Cavallermaggiore and Marene. Wells are typically equipped with submersible pumps that deliver water to storage tanks or directly to the distribution network. Accurate modeling of well output and pump characteristics is essential for reproducing system behavior.

Well Representation Using Negative Demand

Wells were modeled in EPANET using the negative demand approach, where a junction node is assigned a negative base demand equal to the average well production rate. This method is consistent with standard EPANET practice and allows wells to be represented without requiring explicit source elements (reservoirs with fixed heads).

For wells with variable output controlled by tank levels or time schedules, time-varying negative demands were implemented using EPANET patterns linked to control rules.

Pump Curve Development

Pump characteristic curves, which describe the relationship between flow rate and head gain, were derived from manufacturer data sheets and validated using SCADA measurements. The pump curve is typically represented as a polynomial or three-point curve:

$$H = H_0 - rQ^n \quad (4.2)$$

where H is the head gain, H_0 is the shutoff head (head at zero flow), Q is the flow rate, and r and n are curve coefficients.

Manufacturer data provided nominal operating points (flow and head) at design conditions. These points were used to fit the pump curve in EPANET. Where available, SCADA measurements from flow meters and pressure sensors at pump stations were used to validate the curve under real operating conditions.

4.1.6 Control Logic Implementation

Operational control logic governs the behavior of pumps, valves, and other active elements in response to system conditions. The control rules implemented in EPANET replicate the real operational strategies used by Alpi Acque to manage tank levels, maintain pressures, and optimize energy consumption.

Tank Level-Based Pump Controls

The most common control strategy involves activating or deactivating pumps based on tank water levels. A typical rule set includes:

- Pump startup when tank level falls below a minimum threshold (e.g., 20% of capacity)
- Pump shutdown when tank level reaches a maximum threshold (e.g., 80% of capacity)
- Hysteresis margins to prevent rapid cycling (pump restarts only after level drops sufficiently below the maximum threshold)

In EPANET, these controls were implemented using pump curve, [CONTROLS] or [RULES] sections. An example control rule syntax is:

RULE 1

IF TANK Tank1 LEVEL BELOW 1.2
THEN PUMP Pump1 STATUS IS OPEN

RULE 2

IF TANK Tank1 LEVEL ABOVE 2.8
THEN PUMP Pump1 STATUS IS CLOSED

Time-Based Controls

Some pumps operate according to fixed time schedules to take advantage of off-peak electricity rates or to align with anticipated demand patterns. Time-based controls were implemented using EPANET's `AT CLOCKTIME` syntax:

RULE 3

IF SYSTEM CLOCKTIME >= 6:00 AM
AND SYSTEM CLOCKTIME < 10:00 PM
THEN PUMP Pump2 STATUS IS OPEN

Validation of Control Logic

Control logic was validated by comparing simulated pump status logs and tank level trajectories against SCADA operational records. Discrepancies were investigated and resolved by refining threshold values, adjusting hysteresis margins, or correcting misinterpretations of operational procedures. Close collaboration with utility staff was essential for ensuring that the implemented control logic accurately reflected real operational practice.

4.2 Dynamic Stage: Extended-Period Simulation

The Dynamic stage is executed at the beginning of every day, with the data collected by the SCADA platform the day before. Flow measurements at the inlet of every district are used to estimate the dynamic behavior of the nodes inside the district. This dynamic behavior is combined with the base demand of each node obtained in the static stage, as shown in Figure 4.2

The dynamic behavior is applied to each node in EPANET using the `[PATTERNS]` section. Each pattern consists of a sequence of multipliers that scale the node's base demand according to the simulation time. During extended-period simulations (EPS), EPANET applies these multipliers to reproduce realistic daily consumption cycles over multi-day horizons. Figure 4.3 shows an example of such pattern.

Consistency between the spatial component (billing-based base demands) and the temporal component (SCADA-derived time patterns) was assessed by comparing the simulated total consumption against district-level flow meter readings.

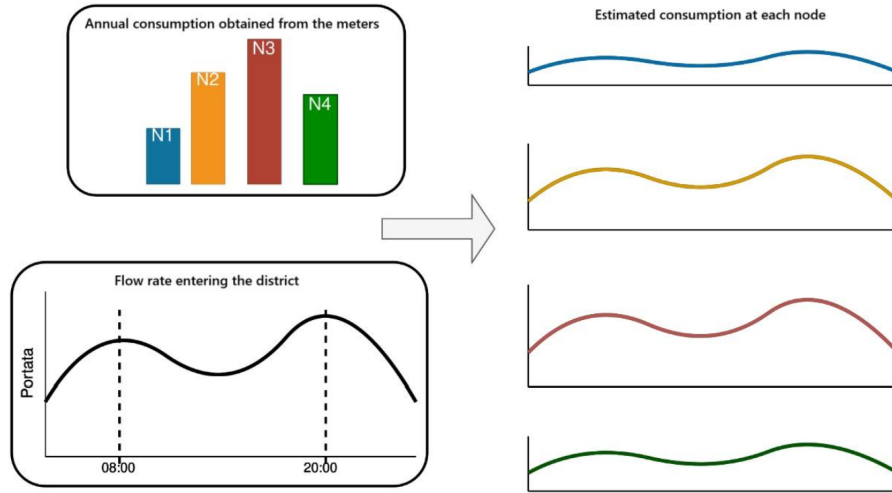


Figure 4.2: Static and dynamic data to estimate nodes demand

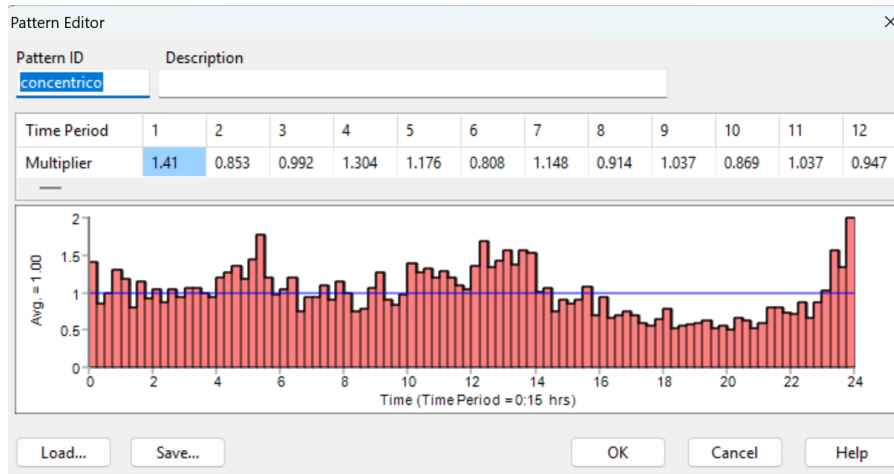


Figure 4.3: Time Pattern example, concentrico district.

When discrepancies were detected, iterative adjustments were made to improve agreement while maintaining the physical interpretation of both data sources.

Extended-Period Simulation (EPS) is a key feature of EPANET that allows the analysis of time-varying hydraulic behavior over multi-day or multi-week periods. Unlike steady-state snapshot simulations, EPS accounts for tank level changes, time-varying demands, and dynamic control actions, providing a realistic representation of network operation.

4.2.1 Output Extraction and Analysis

EPANET generates time-series outputs for all network elements, including nodal pressures, pipe flows, tank levels, and pump status. These outputs were extracted and compared against SCADA measurements to assess the accuracy of the model.

As shown in Figures 4.4 and 4.5 extended-period simulation captures the daily cycling behavior of tanks, the phased activation of pressure and flow variations across the network.

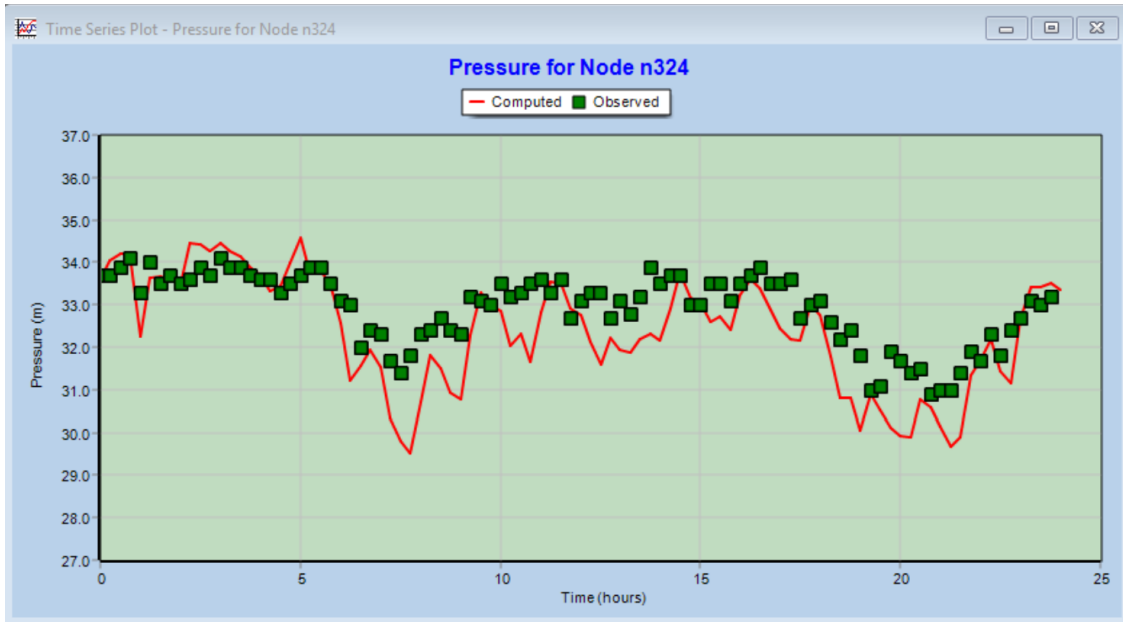


Figure 4.4: An Example of pressure calibration, node n324

4.3 Model Calibration

Model calibration is the process of adjusting uncertain parameters to minimize discrepancies between simulated and observed hydraulic data. Calibration was performed using data sources that were not used for model building:

- Pressure meters at tank outlets and selected network nodes
- Tank level sensors providing continuous water level records
- multi-day pressure measurement at specific nodes from campaigns performed by the utility company.

Using EPANET's calibration data option, plots of the simulated and measured data were produced. Visual inspection allowed to recognize potential issues. Manual

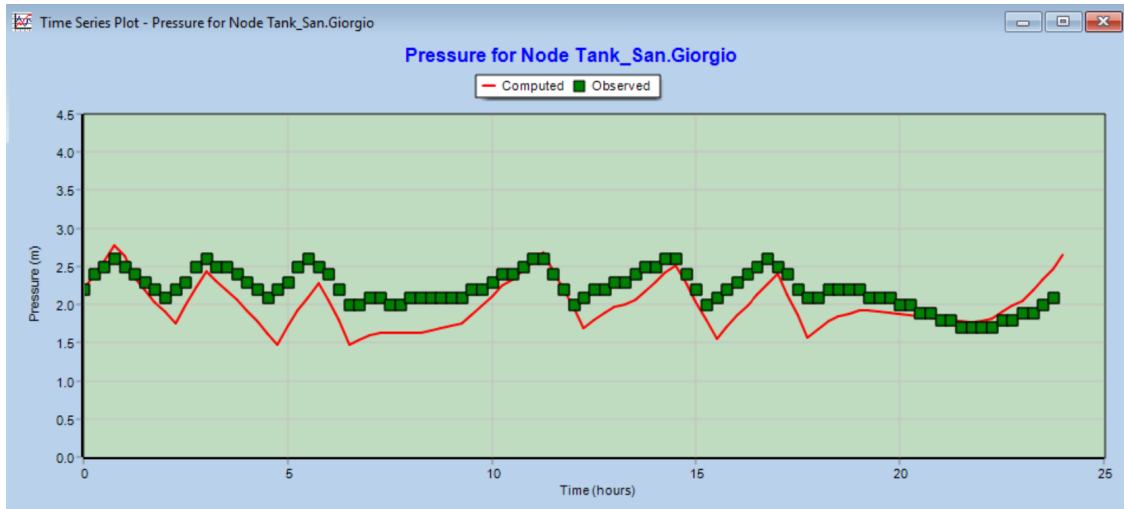


Figure 4.5: An Example of tank level calibration, San Giorgio Tank

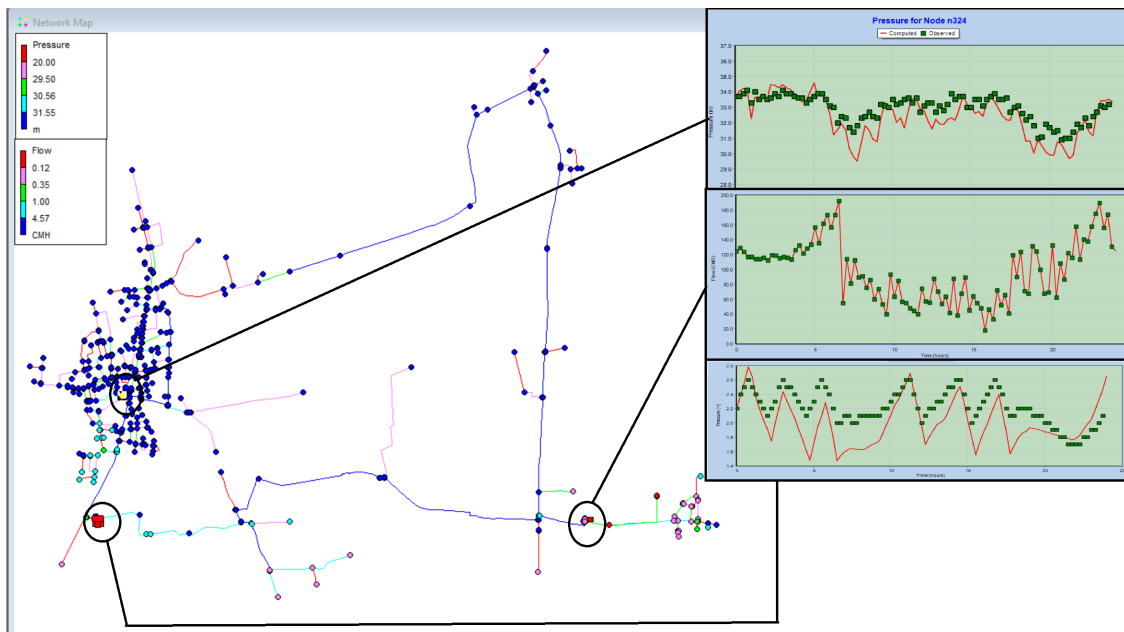


Figure 4.6: Cavallermaggiore Calibration, examples for tank level, flow and pressure sensor

verification of the model was performed to understand the origin of the discrepancy, and the model was fixed accordingly, when possible. In some cases instead, the reason for the differences between simulation and measurements was determined to be a fault in one of the sensors. The utility company was notified and after a check, the sensor error was confirmed. Figures 4.4, 4.5 and 4.6 show some examples of callibration plots.

- **Leakage detection sensitivity:** Sizing districts to provide detectable anomalies relative to measurement noise

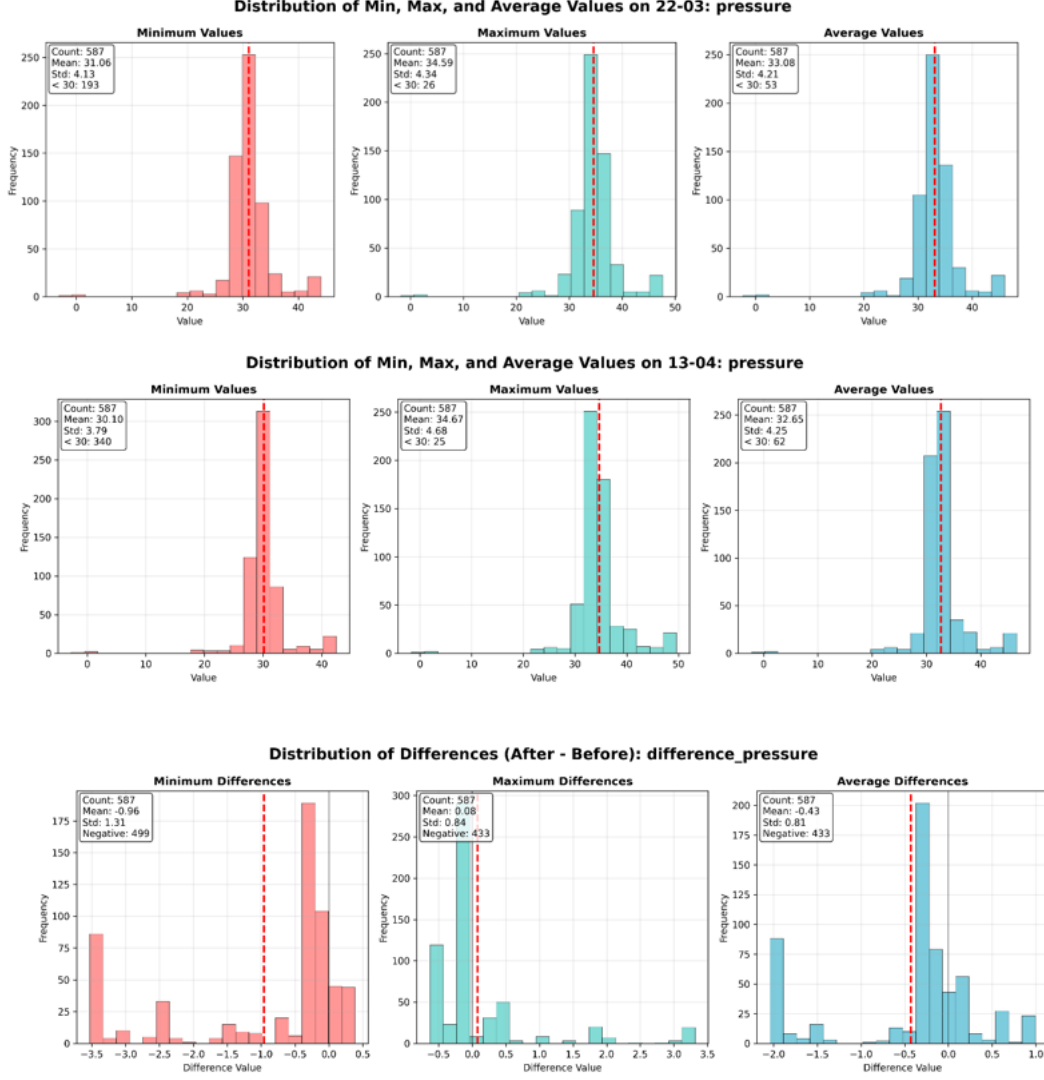


Figure 4.8: Histogram of pressure differentials before and after DMAs

4.4.2 Simulation-Based Evaluation

The hydraulic model was used to simulate the impact of proposed DMA boundaries on pressure distributions and flow patterns.

Using the simulation, an analysis of the drop in pressure was performed, as can be seen in Figure 4.8. This analysis was useful to identify if there were any nodes with a significant pressure drop, or any nodes where the pressure would drop below minimum required value to guarantee quality-of-service to the users.

The DMA boundaries were also implemented in the field by closing selected boundary valves, then as a validation, pressure and flow trending were compared using SCADA. As shown in Figure 4.9, the model is still able to predict the measurements after the DMAs configuration.

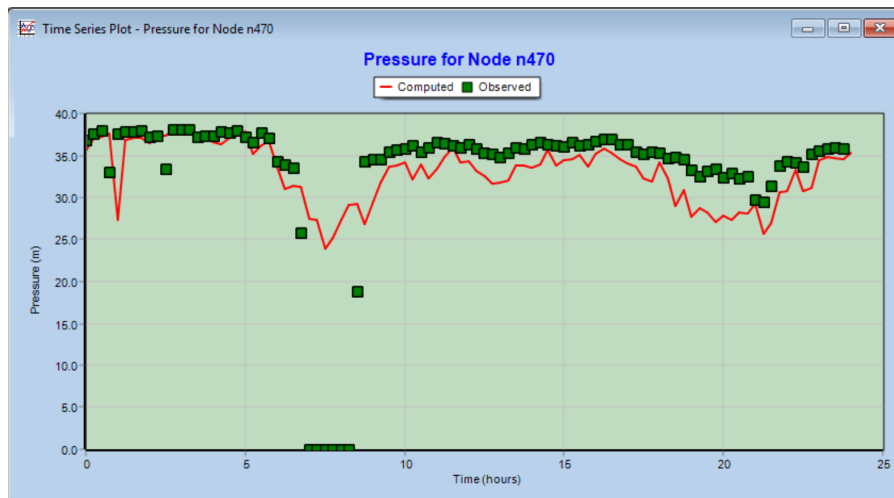


Figure 4.9: Pressure calibration after (DMA) definition for Cavallermaggiore, in node n470

Chapter 5

Part II: Leakage monitoring algorithms

This chapter introduces two methods for leakage monitoring: SCADA-based flow anomaly detection and pressure-based detection/localization using the LILA algorithm with synthetic scenarios

5.1 Flow based leakage detection

5.1.1 Overview and Motivation

Leakage monitoring in water distribution networks requires methodologies that can effectively detect anomalies and localize their sources under practical operational constraints. Building upon the calibrated hydraulic models developed in Chapter 4, this chapter presents two complementary approaches for leakage monitoring in the Marene and Cavallermaggiore networks: data-driven flow-based anomaly detection using machine learning techniques applied to SCADA measurements (Section 5.1), and model-driven pressure-based detection and localization using synthetic scenario generation combined with the LILA algorithm (Section 5.2).

The flow-based approach leverages existing flow measurements at district inlets to identify deviations from normal consumption patterns using the DARTS time-series forecasting framework [35]. While this method provides effective district-level detection using readily available data, it cannot pinpoint leak locations within the district. To address this limitation, a pressure-based methodology is developed utilizing the calibrated EPANET model to generate synthetic leakage scenarios following the BattLeDIM competition framework [2, 8], combined with the LILA (Leakage Identification and Localization Algorithm) [45] for both detection and spatial localization.

LILA was selected for this thesis due to several key advantages: (1) its proven

performance in the BattLeDIM competition, where it successfully detected and localized both abrupt and incipient leaks; (2) its ability to operate effectively with sparse sensor deployments, making it suitable for medium-sized networks with limited instrumentation; (3) its model-based approach that leverages hydraulic simulation to overcome the scarcity of real leak events with concurrent pressure measurements; and (4) its two-stage framework combining statistical process control (CUSUM) for detection with correlation-based pattern matching for localization, providing both robustness and interpretability.

These methodologies provide a practical framework for water utilities with limited sensor infrastructure, enabling both rapid detection of significant leakages and spatial localization to guide field inspection efforts.

5.1.2 DARTS Framework for Time-Series Forecasting and Anomaly Detection

DARTS (Data Analysis with Recurrent Time Series) is a Python library designed for time-series forecasting and anomaly detection [35]. The framework provides a unified interface for training, evaluating, and deploying various forecasting models, including classical statistical methods (ARIMA, exponential smoothing), machine learning approaches (Random Forest, LightGBM), and deep learning architectures (LSTM, Transformer).

For this thesis, DARTS was selected due to its ability to handle time series, integrate exogenous variables, and provide robust anomaly detection capabilities through residual-based scoring. The library’s modular design allows for rapid prototyping and comparison of different modeling approaches, facilitating the selection of the most appropriate method for each district’s consumption characteristics.

Darts first trains a forecasting model on the historical time series in order to learn the expected temporal dynamics. Given an observed series x_t and the model’s predictions \hat{x}_t , the residuals are defined as

$$r_t = x_t - \hat{x}_t.$$

These residuals are then processed through an aggregation step. Using sliding windows of length w , each window is represented as

$$\mathbf{r}^{(t)} = [r_{t-w+1}, \dots, r_t].$$

Feature vectors are extracted from each window and clustered using a k-means clustering aggregator. Let the k-means algorithm learn K cluster centers $\{\mu_k\}$. The anomaly score for window t is computed as the minimum distance to the nearest cluster center:

$$\text{score}_{\text{kmeans}}(t) = \min_k \|f(\mathbf{r}^{(t)}) - \mu_k\|_2.$$

A quantile-based detector is then applied to the distribution of these scores. Given a chosen quantile level α , the threshold is defined as

$$q_\alpha = \text{Quantile}_\alpha(\{\text{score}_{\text{kmeans}}(t)\}).$$

Finally, a time point is flagged as anomalous if

$$\text{score}_{\text{kmeans}}(t) > q_\alpha.$$

In summary, Darts trains the forecasting model, computes residuals, aggregates them through k-means clustering, and applies a quantile-based detector to identify anomalies.

5.1.3 Model Selection and Configuration

Based on preliminary experiments, an LSTM (Long Short-Term Memory) network was selected as the primary forecasting model primarily due to the limited availability of historical consumption data. In scenarios with relatively few data points, LSTM networks have been observed to perform better than alternative machine learning models because of their ability to effectively capture temporal dependencies without requiring extremely large datasets. Their internal memory cells enable the model to learn patterns over time, even when the training data are sparse, and to generalize better than simpler architectures.

For each district, two years of historical consumption data were used, with a temporal resolution of 15 minutes. This corresponds to 96 measurements per day, resulting in a total of approximately 70080 samples per district, which provides a sufficiently rich dataset for the LSTM to learn both short-term fluctuations and long-term patterns.

The chosen model architecture is designed to balance complexity and regularization to prevent overfitting, and includes:

- Input sequences of 96 timesteps, corresponding to 24 hours of consumption data at 15-minute intervals
- 4 stacked LSTM layers, each with 128 hidden units, allowing the network to model increasingly abstract temporal patterns
- Dropout regularization with a rate of 0.2, applied to each layer to reduce the risk of overfitting
- Training for 20 epochs with a batch size of 32, ensuring adequate learning without excessive computational cost
- Optimization using the Adam optimizer with a learning rate of 0.001, which provides a good balance between convergence speed and stability

To further account for the heterogeneous consumption characteristics across the network, the model was trained independently for each district. This approach allows the LSTM to adapt to local consumption behaviors, capturing both short-term fluctuations and long-term trends, while maintaining robustness in situations with limited historical data.

5.2 Pressure based leakage detection and localization

Pressure-based leakage localization requires relating observed pressure patterns to potential leak locations through hydraulic modeling. However, documented leakage events with concurrent high-quality pressure measurements are scarce in operational networks. To address this challenge, this thesis adapts the strategy followed by the BattLeDIM (Battle of the Leakage Detection and Isolation Methods) competition [2, 8], which established a benchmark framework for evaluating detection and localization algorithms using synthetic data.

The BattLeDIM approach involves generating a network model representing the real system and a "noisy" model with uncertain parameters to simulate imperfect knowledge conditions. Synthetic SCADA measurements are generated from simulations, including pressure and flow time series with realistic characteristics and multiple leakage scenarios with varying sizes, locations, and temporal patterns. While BattLeDIM used an idealized synthetic network (L-TOWN), this thesis applies the methodology to the real Cavallermaggiore network using the calibrated EPANET model developed in Chapter 4.

For detection and localization, the LILA (Leakage Identification and Localization Algorithm) [45], validated in the BattLeDIM competition, is adapted to the Cavallermaggiore case. LILA employs a two-stage framework: first, pairwise linear regression models between pressure sensors are constructed during leak-free periods, and deviations from these relationships are monitored using CUSUM (Cumulative Sum) statistical process control charts to detect anomalies; second, once a leak is detected, the Most Affected Sensor (MAS) is identified, and the hydraulic model simulates candidate leak locations to find the best correlation match with observed pressure patterns. This model-based approach enables effective localization even with sparse sensor deployments typical of medium-sized municipal networks.

5.2.1 Simulation of leakages with BattleDIM approach for Cavallermaggiore

The synthetic leakage scenario generation follows the BattLeDIM methodology adapted to the calibrated Cavallermaggiore network model. The process involves

creating both a baseline “clean” model and a “noisy” model to represent realistic uncertainty conditions.

Model Uncertainty Representation:

To replicate the imperfect knowledge conditions, a noisy version of the calibrated model was generated by introducing controlled 10% perturbations to uncertain parameters. The noise was applied systematically to:

- **Pipe roughness coefficients:** $\pm 10\%$ variation around calibrated values to simulate pipe aging uncertainty and calibration errors
- **Pipe lengths and diameters:** $\pm 10\%$ variation to account for survey uncertainties and database inaccuracies
- **Nodal base demands:** $\pm 10\%$ variation to represent demand estimation uncertainty from billing data and consumption pattern variability

Hydraulic Simulator Selection:

For the hydraulic simulation backend, EpanetSimulator was selected over WNTR’s Python-based solver. The EpanetSimulator uses the native EPANET engine and provides significantly better computational performance for extended-period simulations, completing annual simulation cycles in approximately 5-10 minutes compared to 2-5 hours with the WNTR Python solver. This performance advantage is particularly important for generating comprehensive datasets with multiple leakage scenarios spanning full annual cycles.

Reservoir and Pump Configuration:

During the initial leak scenario generation, simulations encountered negative pressure values, particularly when adding significant leakages to the network. Since the simulated network model represented the water supply to the tank using negative base demands at the tank node to simulate inflow, this approach lacked the flexibility to automatically maintain tank levels within operational ranges during leak events, causing the tank level to drop minimum thresholds and resulting in negative pressures in downstream areas.

Following EPANET modeling practices, this issue was resolved by introducing a reservoir (Reservoir_SanGiorgio) with fixed head at 310 m elevation, representing an infinite water source such as an aqueduct connection. A pump (Tank_Pump) was added to transfer water from the reservoir to the tank, controlled by simple rules based on tank levels:

$$\text{Pump} = \begin{cases} \text{OPEN} & \text{if Tank level} < 1.3 \text{ m} \\ \text{CLOSED} & \text{if Tank level} > 2.7 \text{ m} \end{cases} \quad (5.1)$$

This configuration ensures that the tank automatically refills when levels drop due to increased demand or leakage, maintaining realistic pressure conditions throughout the network and preventing the occurrence of physically impossible negative pressures during extended simulations.

5.2.2 LILA for detection

Leakage detection in water distribution networks relies on analyzing the pressures measured at different nodes and identifying deviations from their expected behaviour. To characterise these expected relations, a simple linear model is built for each node i , using the pressures of the neighbouring nodes as explanatory variables.

In its simplest form, the expected pressure at node i can be approximated as:

$$\hat{P}_i(t) = a_{0,i} + \sum_{j \in \mathcal{N}(i)} a_{j,i} P_j(t),$$

where $P_j(t)$ denotes the pressure at a neighbouring node j , $\mathcal{N}(i)$ is the set of nodes hydraulically connected to i , and $a_{0,i}, a_{j,i}$ are coefficients estimated from leak-free data.

The deviation between the measured and the predicted pressure,

$$\text{MRE}_i(t) = P_i(t) - \hat{P}_i(t),$$

represents the residual error of the model. Under normal operating conditions, $\text{MRE}_i(t)$ fluctuates around zero. When a leak occurs, the additional unmodelled water loss perturbs the hydraulic equilibrium and produces systematic deviations in the residuals, especially at nodes located close to the leakage.

By comparing the predicted pressures with the observed measurements, it is possible to compute the Model Residual Error (MRE) for each node. Under normal, leak-free conditions, the MRE fluctuates around zero, reflecting only measurement noise or minor natural variations. When a leakage occurs, the additional unaccounted flow produces systematic deviations in the MRE, particularly at nodes located near the leak.

The linear regression-based MRE allows the identification of the most affected sensors at any given time, i.e., nodes with the largest deviations from their predicted pressures. The magnitude of these deviations can be aggregated into a score for each node, which provides a ranking of nodes according to their likelihood of being impacted by a leak. By monitoring the time evolution of these scores, it is possible to detect the start of a leakage event and evaluate its impact across the network.

This methodology works for both single and multiple simultaneous leaks. In the case of multiple leaks, the aggregated MRE highlights all critical areas, although disaggregation into individual leaks may require further analysis. Overall, the combination of linear regression modeling and residual error analysis provides

a systematic, data-driven framework for leakage identification, enabling operators to exploit the existing network of pressure sensors effectively, even when the precise location of leaks is unknown.

This approach is inspired by the LILA (Leak Identification in Large-scale water networks) methodology [45], which applies a similar combination of linear regression and residual analysis to identify critical nodes and detect leakages, providing a validated framework for data-driven leak detection in complex water distribution systems.

5.2.3 LILA for localization

In this work, the leak localization methodology originally proposed by LILA was adapted to our network and simulation setup, since the complete LILA code is not publicly available. The underlying idea remains the same: to characterize how pressure at monitoring nodes responds to leaks occurring at different pipe locations, allowing the identification of likely leak positions.

In a preliminary stage, a synthetic leak is simulated on every pipe of the network. For each simulated leak, the induced pressure variation ΔP at all pressure sensors is computed. This process yields a dataset in which each pipe is associated with a vector of pressure deviations—its hydraulic signature. These signatures describe how a leak at that specific pipe propagates through the system and affects the monitored nodes.

Using this dataset, it is possible to determine, for each sensor, which pipes have the strongest influence on its readings. Since each simulation quantifies the magnitude of the pressure drop at each sensor, the pipes can be ranked according to their ΔP values for a given node. For each monitoring node, the 30 most influential pipes are selected, representing the most hydraulically plausible leak candidates from that sensor’s perspective.

To further refine the candidate set and avoid overly broad localization, an additional filtering step is applied. Specifically, for each sensor we retain only the pipes for which that sensor is the most affected node in the corresponding simulation, i.e., the node exhibiting the maximum ΔP . This ensures that the remaining pipes are both highly influential and directly connected to the local pressure response observed at the sensor.

Through this two-step selection and filtering process, each sensor is associated with a compact and hydraulically consistent set of candidate pipes. The resulting framework significantly reduces the spatial search area and enhances the precision of the leak localization procedure, providing a reliable basis for subsequent diagnostic or operational actions.

Chapter 6

Results from Leakage monitoring strategies

Here we present the results obtained from the two proposed leakage monitoring strategies, covering both flow-based anomaly detection and pressure-based detection and localization.

6.1 Leakage detection based on flow sensors

A portion of the experiments in this work concerns anomaly detection using flow sensors. In fact, we initially worked in the city of Marene, where only flow sensors are available; for this reason, anomaly detection was carried out based on them. Specifically, each flow sensor was associated with a district of the city, and in turn each district was assigned its own anomaly-detection model, given that every district has a different data distribution.

Before applying the methodology to Marene, the leak detection approach was first tested on the *BattleDim* benchmark dataset (2018). The performance metrics obtained in this preliminary evaluation are reported in Table 6.1. These metrics are computed as follows:

- **Precision:** $\text{Precision} = \frac{TP}{TP+FP}$
- **Recall:** $\text{Recall} = \frac{TP}{TP+FN}$
- **F1-score:** $F1 = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$
- **AUC:** Area under the ROC curve, measuring overall discrimination ability.

Figure 6.1 shows the comparison between detected anomalies and the actual anomalies on the *BattleDim* dataset. This illustrates how well the model identifies anomalous flow patterns.

Metric	Value
Precision	0.745
Recall	0.84
F1-score	0.79
AUC	0.84

Table 6.1: Performance metrics of the leak detection model on the BattleDim dataset 2018.

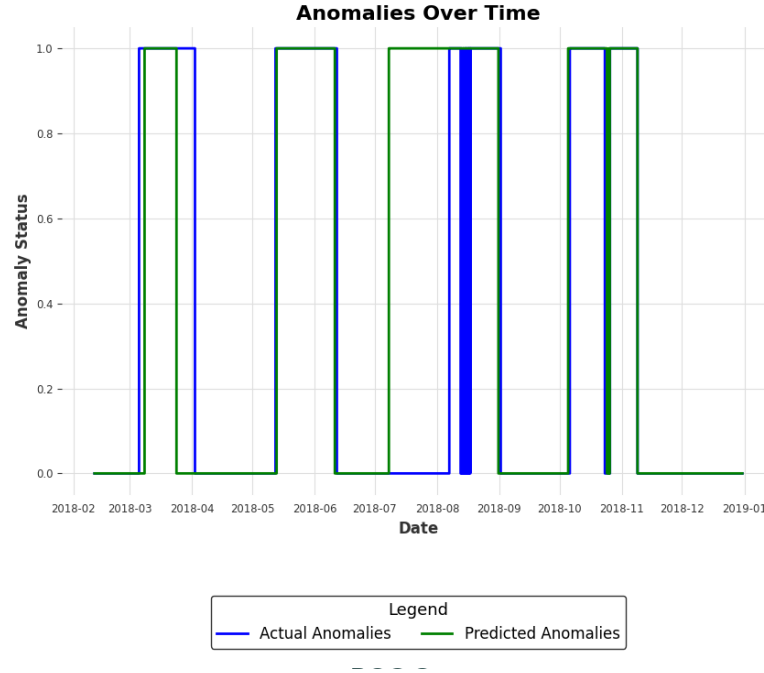


Figure 6.1: Detected anomalies vs actual anomalies on the BattleDim dataset.

Figure 6.2 shows the receiver operating characteristic (ROC) curve and the corresponding AUC, highlighting the model's ability to discriminate between normal and anomalous flow conditions.

Finally, Figure 6.3 shows the confusion matrix for the model predictions, providing a detailed view of true positives (TP), false positives (FP), true negatives (TN), and false negatives (FN).

Once the leak detection methodology demonstrated satisfactory performance on the BattleDim dataset, it was applied to all districts of the city of Marene. However, during this phase, several challenges emerged. Specifically, not all leaks were properly reported in the historical dataset, and some detected leaks were signaled with delays of several days. This made it difficult to define quantitative

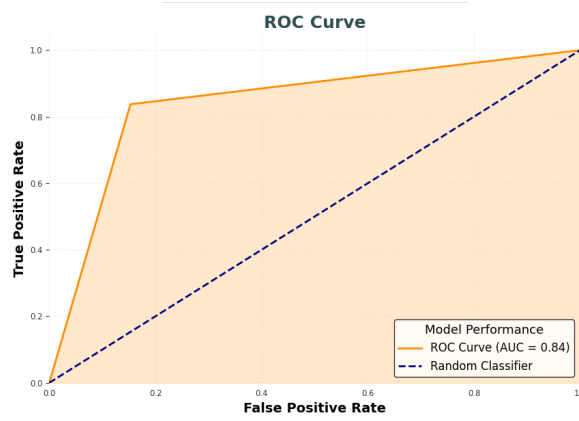


Figure 6.2: ROC curve for the leak detection model on the BattleDim dataset. The AUC value is 0.84, indicating strong classification performance.

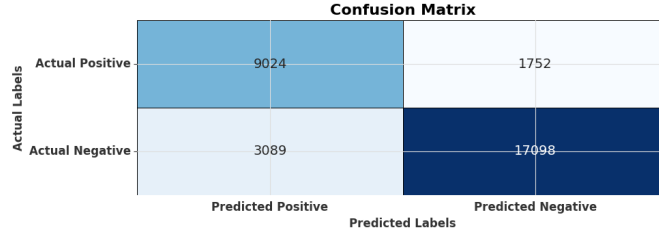


Figure 6.3: Confusion matrix of the leak detection model on the BattleDim dataset 2018 year.

performance metrics for the Marene case, as a direct comparison with ground truth was not always possible (see Figure 6.4).

Despite these limitations, qualitative analysis of the plots shows that, for example, in the district of Pellaverne, the model is able to signal anomalies with reasonable timing (Figure 6.5). This indicates that the approach can still provide valuable operational insights, allowing early detection of abnormal flow patterns and supporting proactive maintenance actions.

As can be observed from Figures 6.7 and 6.6, the behavior of the leak detection model in Marene varies across districts. Notably, the sensors were not always fully operational or correctly configured. For instance, in the district of Marconi, a significant number of false positives were recorded. This was primarily due to misconfigurations or temporary malfunctions of the flow sensors, which prevented the model from running correctly on that data.

Such issues underscore a key practical consideration when deploying anomaly detection models in real-world water networks: unlike benchmark datasets such as BattleDim, where the data is clean and well-structured, real sensors can exhibit missing readings, noise, or misconfigurations. These factors can significantly affect

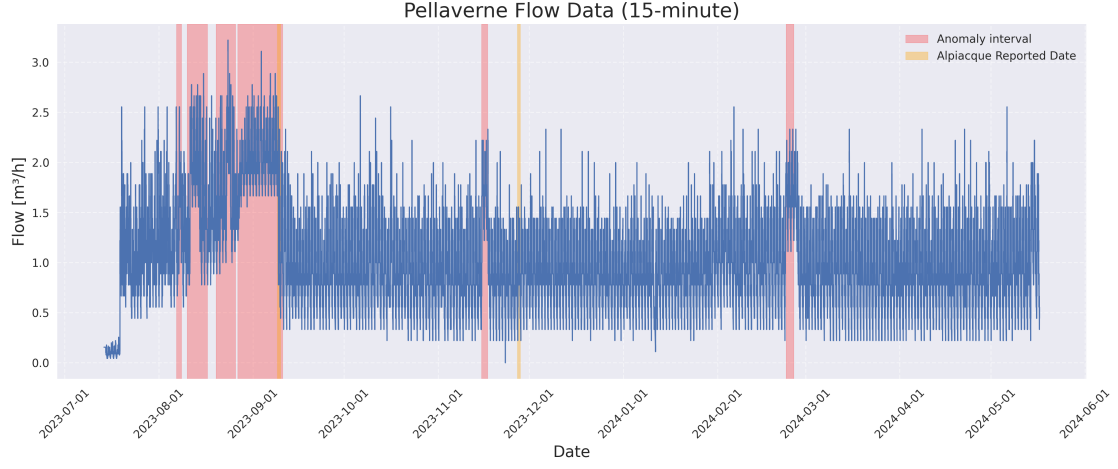


Figure 6.4: Overview of leak detection results for the district of Pellaverne, Marene, using 15-minute intervals. The variability in reporting times and occasional missing detections across all districts highlights the challenges in computing standard performance metrics.

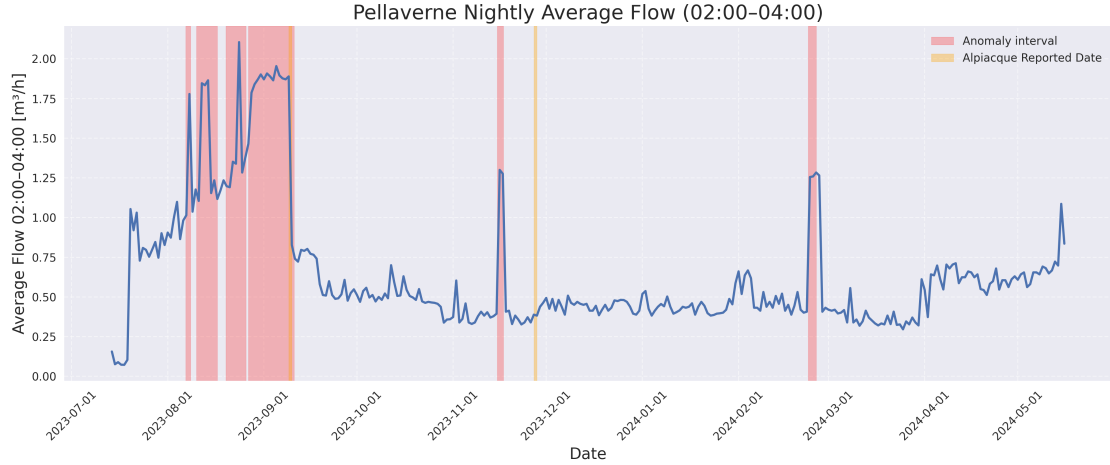


Figure 6.5: Detected anomalies in the district of Pellaverne, Marene, for the nightly period between 2:00 and 4:00 a.m. These hours provide a clearer view of the flow behavior, allowing better identification of potential leaks. The model is able to signal anomalies with reasonable timing, supporting timely intervention.

the reliability of anomaly detection, making it crucial to account for sensor performance and maintenance when interpreting results and planning interventions.

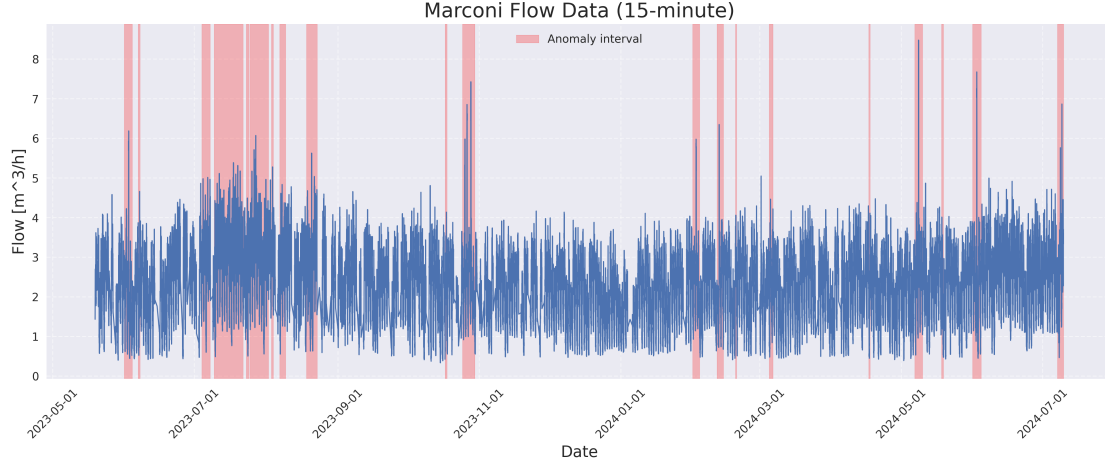


Figure 6.6: Overview of leak detection results for the district of Marconi, Marene, using 15-minute intervals. The variability in reporting times, false positives, and occasional missing detections highlights the practical challenges of deploying real-world sensors and explains why computing standard performance metrics was not feasible.

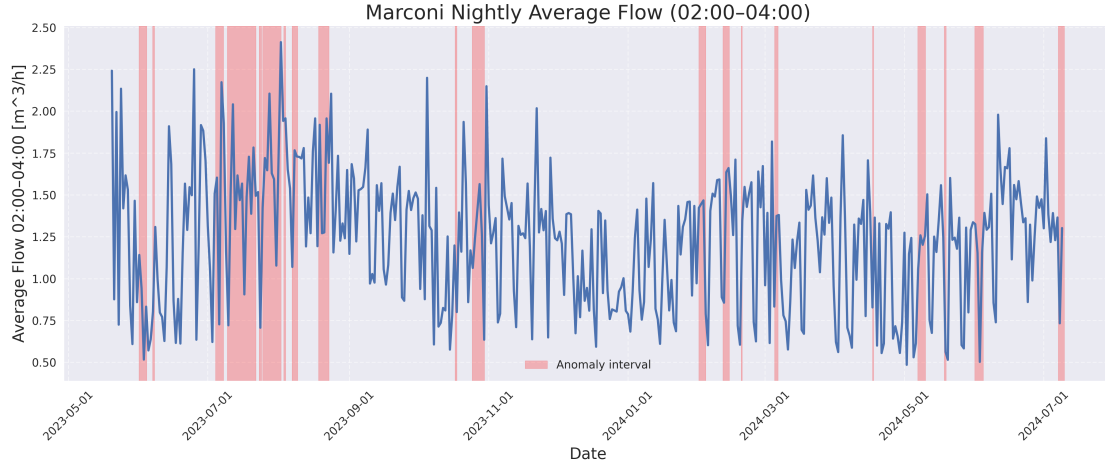


Figure 6.7: Detected anomalies in the district of Marconi, Marene, during the nightly period between 2:00 and 4:00 a.m. Compared to Pellaverne the signal has more peaks, due the misconfiguration problems

6.2 Leakage detection and localization based on pressure sensors

6.2.1 Initial evaluation on BattleDIM

Since the LILA algorithm was presented during the BattleDIM competition, there was no need for evaluation against this dataset. However, only the detection component of the algorithm has been released with clear usage instructions. For this reason, we re-implemented the localization component following the description of the algorithm in the paper, and first tested it on the BattleDIM dataset for District A. The purpose of this preliminary evaluation was to verify that our implementation is consistent with the original methodology.

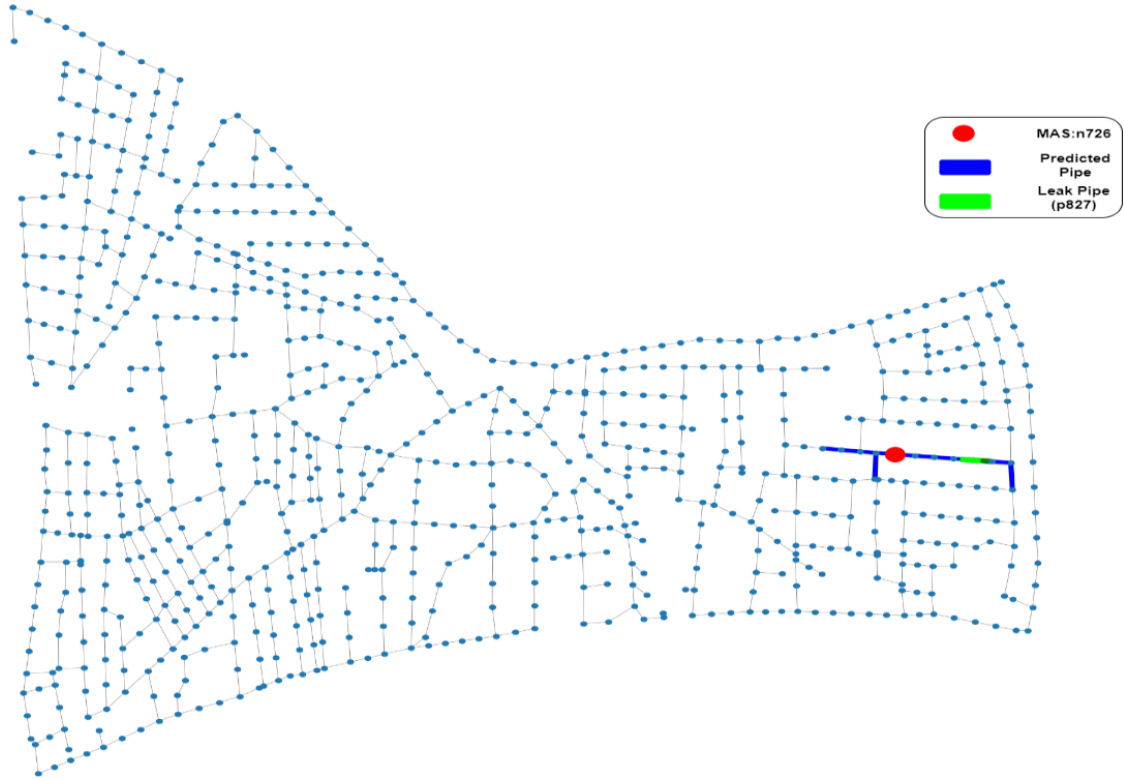


Figure 6.8: Leak localization results for the BattleDIM 2019 dataset for MAS 276. Plot shows the pressure response of an individual sensor to simulated leaks. The model is able to rank candidate pipes according to the likelihood of being the leak location, providing a clear indication of the most probable leak positions.

The results, shown in Figures 6.8, indicate that the localization model is able to provide, with reasonable accuracy, an ordered list of candidate pipes. Each pipe is ranked according to the likelihood of being the leak location, demonstrating that

the reimplemented algorithm effectively reproduces the leak localization behavior reported in the original LILA study.

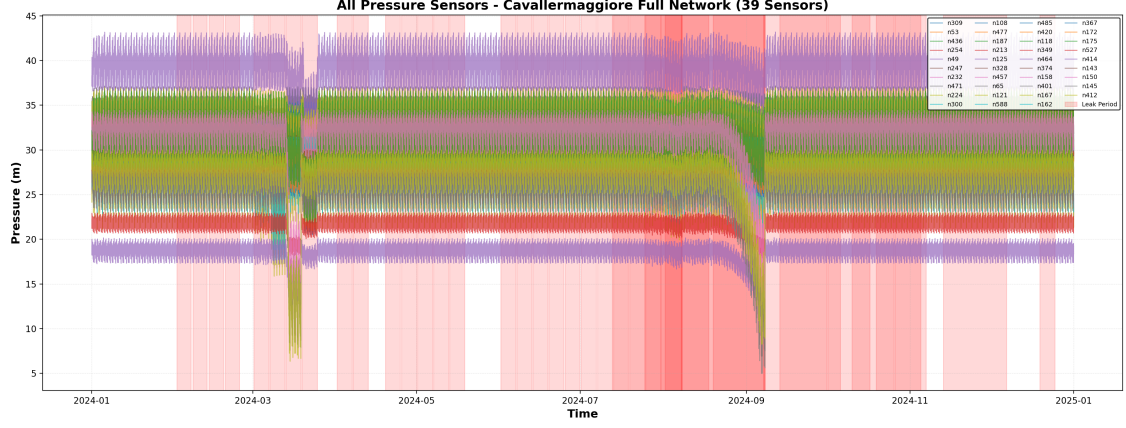


Figure 6.9: All pressure sensors together

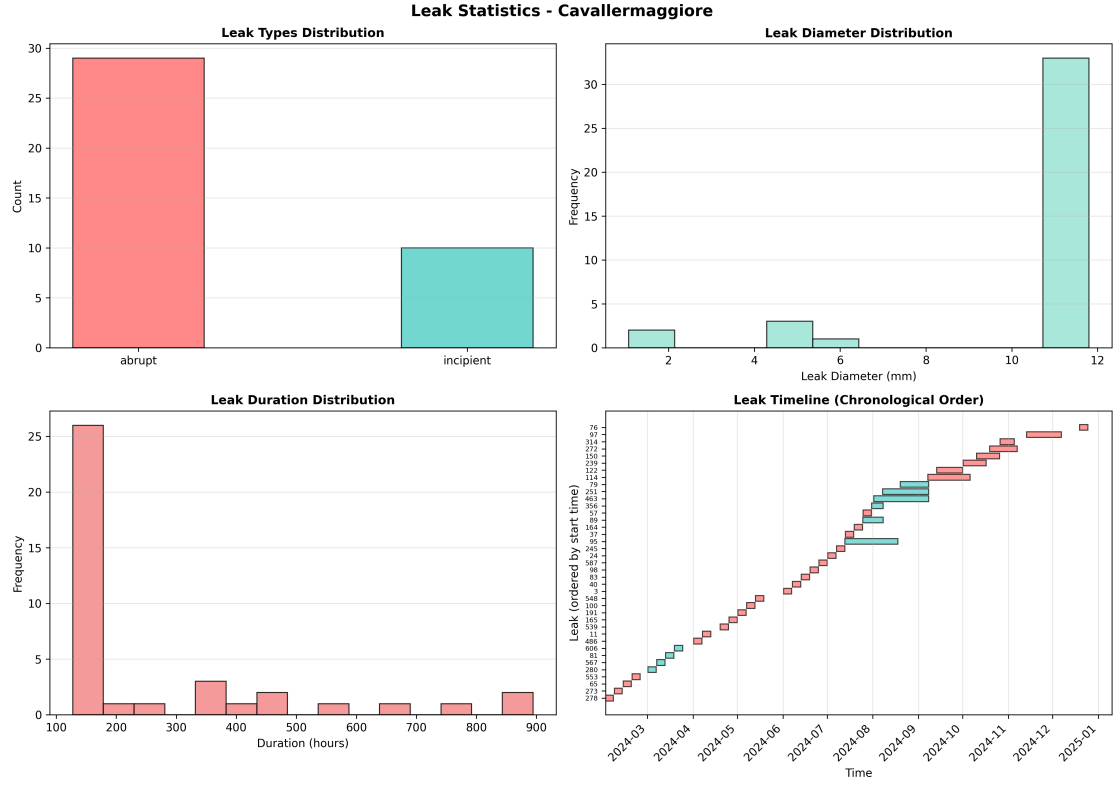


Figure 6.10: Statistics of simulated leakages

6.2.2 Data generation

For the Cavallermaggiore network, 39 leakages scenarios were added. The simulation is run for a full year (2024), where the first month has no leakages and is used for model training, while the rest of the years contains the 39 leakages. A total of 39 pressure sensors were simulated, distributed across all the districts. Figure 6.9 shows the data from the virtual sensors overlayed with the time periods with leakages. Figure 6.10 summarizes some important statistics of the leakages.

6.2.3 Leakage identification step

To evaluate the methodology described on previous chapter, the city of Cavallermaggiore was considered as a case study. The network was divided into multiple districts, and for each district, the linear regression-based leakage detection algorithm was executed independently. This district-wise approach accounts for the different hydraulic characteristics and operational conditions across the network, allowing a more precise identification of critical nodes.

For each district, the algorithm computes the Most Affected Sensors (MAS) by evaluating the magnitude of the Model Residual Error (MRE) over time. The MAS indicate the nodes that are likely to be influenced most strongly by potential leaks, providing a ranked list that can guide further investigation or intervention.

Figures 6.11–6.17 show the MAS plots for all seven districts of Cavallermaggiore. In each plot, the nodes with the highest deviations from their predicted pressures are highlighted, illustrating how the algorithm captures the spatial distribution of anomalies and identifies the most critical monitoring points in the network.

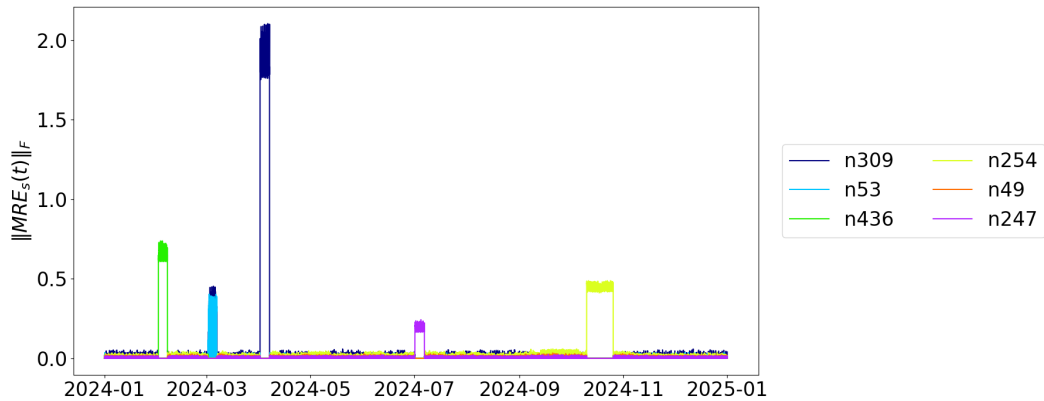


Figure 6.11: Most Affected Sensors (MAS) for District Concentrico of Cavallermaggiore. Nodes with larger MRE values indicate higher likelihood of being influenced by a leak.

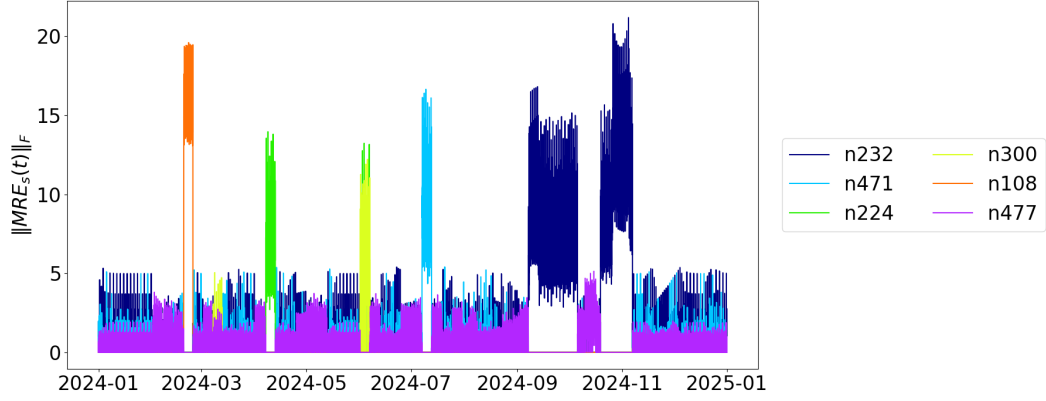


Figure 6.12: Most Affected Sensors (MAS) for District Via Bra of Cavallermaggiore.

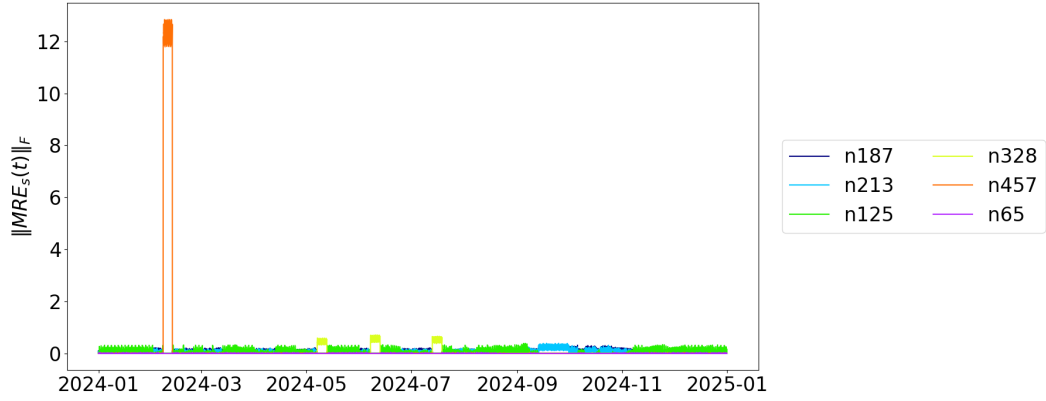


Figure 6.13: Most Affected Sensors (MAS) for District Via Europa of Cavallermaggiore.

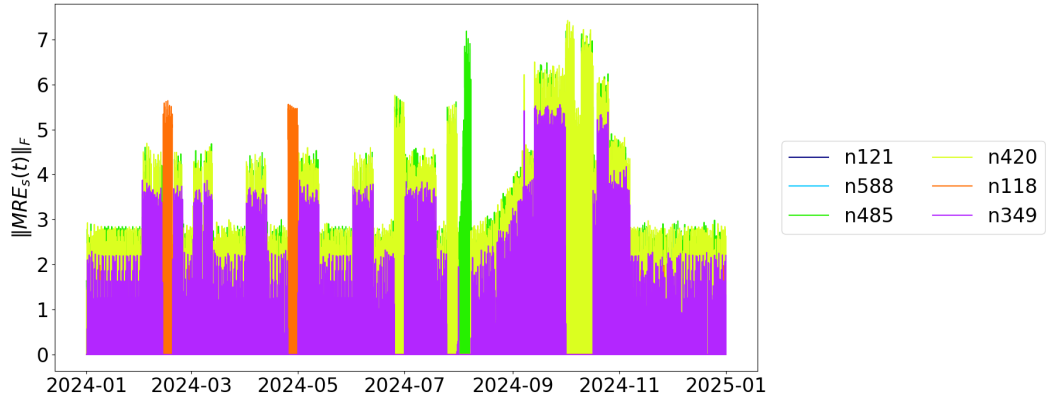


Figure 6.14: Most Affected Sensors (MAS) for District Via Roma of Cavallermaggiore.

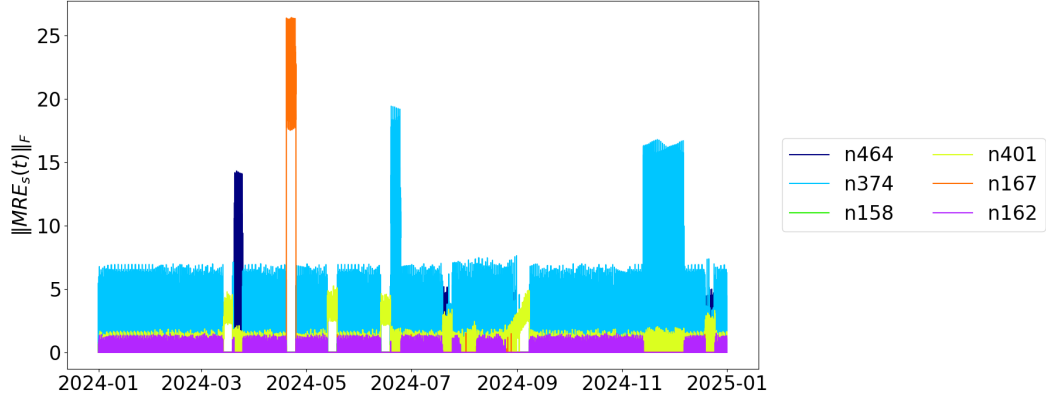


Figure 6.15: Most Affected Sensors (MAS) for District Provinciale of Cavallermaggiore.

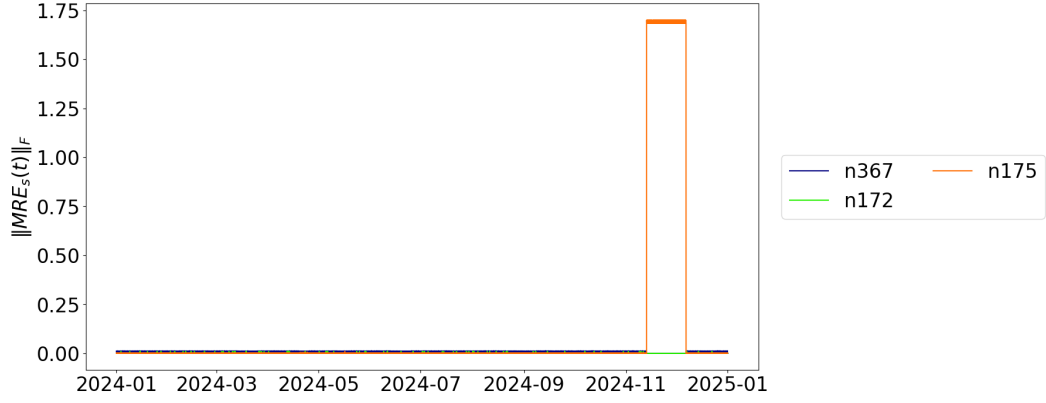


Figure 6.16: Most Affected Sensors (MAS) for District Foresto of Cavallermaggiore.

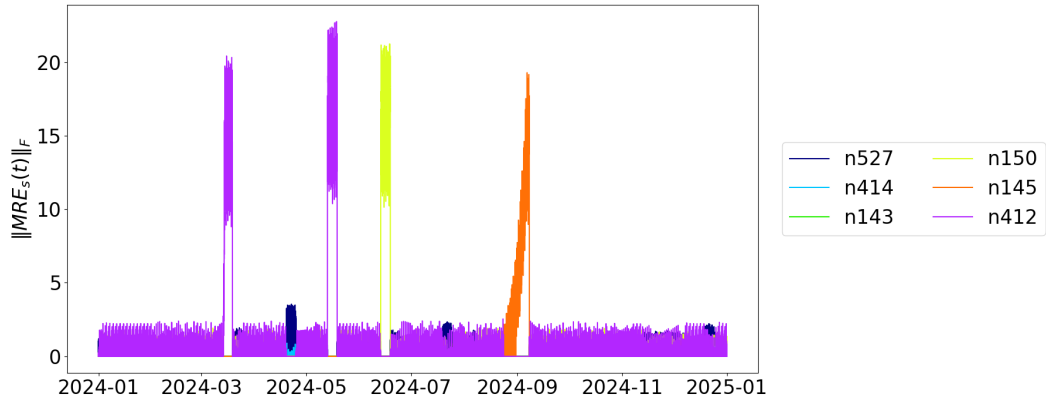


Figure 6.17: Most Affected Sensors (MAS) for District Madonna Del Pilone of Cavallermaggiore.

6.2.4 Leakage localization step

The results of the adapted leak localization methodology are illustrated for the Concentrico district of Cavallermaggiore. The analysis focuses on the identification of likely leak positions using the MAS (Most Affected Sensors) approach, comparing predicted leaks with the ground truth events. Figure 6.18 shows the network layout of Cavallermaggiore, highlighting the predicted leak pipe, the actual leak pipe (ground truth), and the sensor identified as the MAS for this event. This visualization provides an immediate understanding of how closely the localization algorithm can approximate the true leak location and which sensors are most informative for the detection.

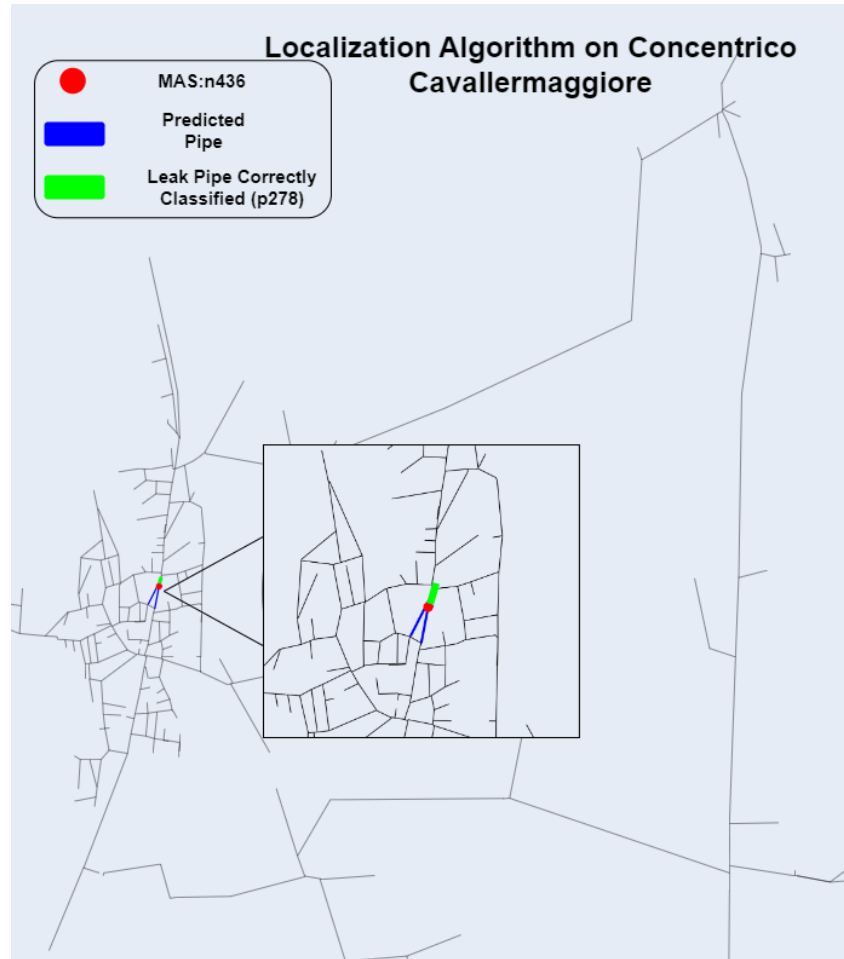


Figure 6.18: Leak localization results in the *Concentrico* district of Cavallermaggiore. The red marker indicates the Most Affected Sensor (MAS), the blue line shows the predicted leak pipe, and the green line corresponds to the ground truth leak pipe.

As can be observed, the algorithm correctly identifies the pipe associated with the leak, and the MAS sensor is positioned in close hydraulic proximity to the leak. This confirms that the adapted methodology effectively reduces the spatial search area and provides actionable insights for network monitoring and timely intervention. Overall, this case demonstrates that even with partial access to the original LILA implementation, the adapted procedure is able to localize leaks accurately within the network, highlighting both the candidate pipes and the key sensors involved in detecting pressure deviations.

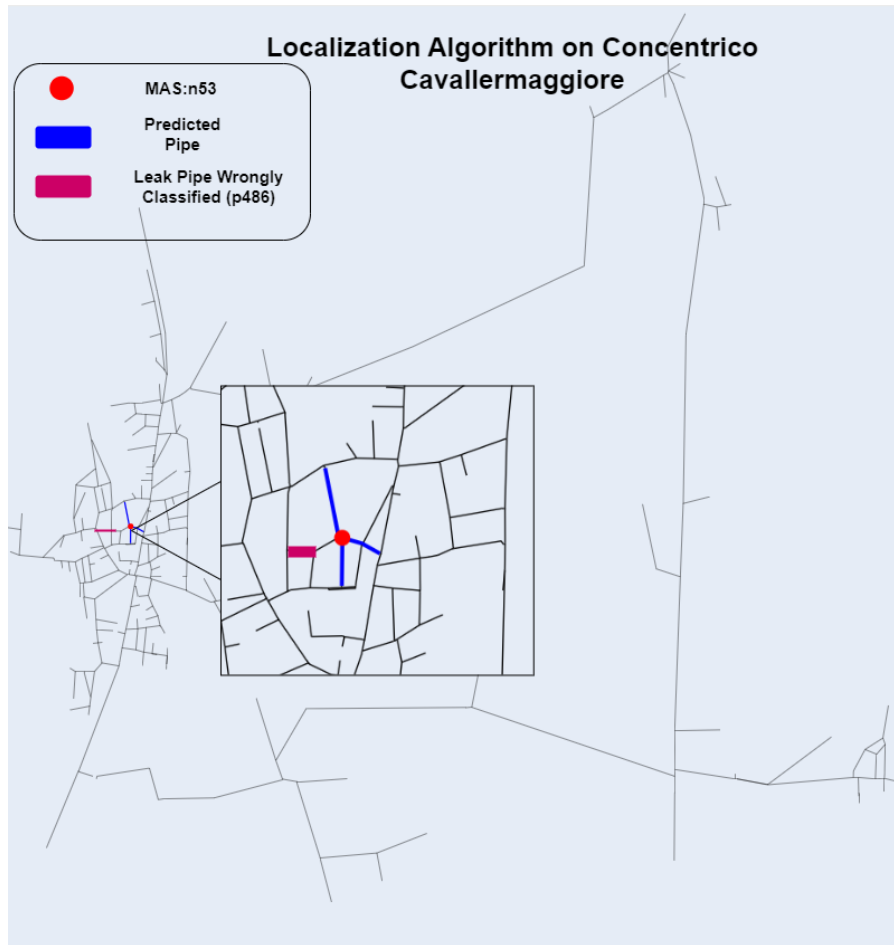


Figure 6.19: Leak localization results for the misclassified cases in the *Concentrico* district of Cavallermaggiore. Although the predicted leak pipes (in blue) do not exactly match the ground truth pipes (in magenta), the algorithm still identifies locations in close hydraulic proximity to the actual leaks. The red marker indicates the Most Affected Sensor (MAS).

However, it is important to note that not all cases are perfectly identified by the adapted methodology. As illustrated in Figure 6.19, there are situations where the

algorithm does not pinpoint the exact leaking pipe but instead predicts a neighbouring pipe with similar hydraulic influence. In the visualization, these incorrectly classified pipes are highlighted in magenta, emphasizing the areas where the model's prediction deviates from the ground truth. Despite these inaccuracies, the algorithm consistently identifies pipes that are hydraulically close to the actual leak location. This proximity still provides valuable operational guidance, substantially narrowing the search area and allowing field operators to concentrate their inspections where pressure anomalies are most likely to originate. Furthermore, the MAS sensor identified in these cases continues to serve as a reliable indicator of the zone most affected by the leak event.

6.2.5 Results summary and discussion

The results reported in Table 6.2 highlight substantial differences in leak localization performance across the various districts of the Cavallermaggiore network. Among all areas, the *Concentrico* (CON) district shows the most reliable behaviour. In this district, all detectable leaks are successfully identified and correctly localized within the ranked candidate list. This strong performance is largely attributable to the specific sensor configuration: since the artificially generated leaks happened to occur near existing pressure sensors, the hydraulic disturbances produced by each event were captured clearly and unambiguously. As a consequence, the MAS identification was straightforward, and the localization algorithm was able to isolate the correct pipe or its immediate neighbours.

In contrast, the *Via Roma* (ROM) district represents the most challenging scenario. Here, several leaks remain undetected or cannot be accurately localized despite being detected. A key reason is the hydraulic role of this district within the network: ROM acts as a primary supply corridor feeding multiple downstream areas. Because of this, the pressure dynamics are strongly influenced by global network behaviour rather than by purely local conditions. When a district is responsible for feeding the rest of the system, pressure sensors within it may undergo significant variations originating from demand fluctuations, valve operations, or distant hydraulic interactions. These background variations can mask or dilute the pressure signatures generated by leaks, making the MAS assignment less reliable and reducing the contrast between normal and faulty conditions.

Moreover, since leak-induced pressure drops propagate across a wider portion of the network, the spatial sensitivity of the sensors becomes less sharp. Instead of capturing a localized disturbance, sensors in ROM may register a more distributed and attenuated pressure response. This results in lower ranking accuracy for the true leaking pipe and, in some cases, prevents the algorithm from including the correct pipe among the top candidates.

Dis	Pipe	Start	End	Size	Type	Det.	MAS	Loc.	Rank
CON	278	02-01	02-06	0.0117	ABR	YES	436	YES	7
CON	280	03-01	03-06	0.0111	INC	YES	309	YES	35
CON	486	04-01	04-06	0.0110	ABR	YES	309	YES	4
CON	191	05-01	05-06	0.0115	ABR	NO	-	-	-
CON	24	07-01	07-06	0.0115	ABR	YES	247	YES	15
CON	463	08-01	09-07	0.0012	INC	NO	-	-	-
CON	150	10-10	10-25	0.0110	ABR	YES	254	YES	1
BRA	553	02-19	02-24	0.0114	ABR	YES	108	NO	-
BRA	567	03-07	03-12	0.0118	INC	NO	-	-	-
BRA	11	04-07	04-12	0.0111	ABR	YES	224	YES	18
BRA	3	06-01	06-06	0.0111	ABR	YES	300	YES	8
BRA	245	07-07	07-12	0.0114	ABR	YES	471	YES	3
BRA	114	09-07	10-05	0.0117	ABR	YES	232	YES	26
BRA	314	10-26	11-04	0.0051	ABR	YES	232	YES	5
EUR	273	02-07	02-12	0.0112	ABR	YES	457	YES	3
EUR	100	05-07	05-12	0.0111	ABR	YES	328	YES	23
EUR	40	06-07	06-12	0.0115	ABR	YES	328	YES	16
EUR	37	07-13	07-18	0.0113	ABR	YES	328	YES	15
EUR	251	08-07	09-07	0.0112	INC	YES	187	NO	-
EUR	122	09-13	09-30	0.0113	ABR	YES	213	YES	5
MDP	81	03-13	03-18	0.0111	INC	YES	412	YES	5
MDP	548	05-13	05-18	0.0115	ABR	YES	412	YES	7
MDP	83	06-13	06-18	0.0110	ABR	YES	150	YES	4
MDP	79	08-19	09-07	0.0115	INC	YES	145	YES	3
PRO	606	03-19	03-24	0.0058	INC	YES	464	YES	29
PRO	539	04-19	04-24	0.0116	ABR	YES	167	YES	16
PRO	98	06-19	06-24	0.0051	ABR	YES	374	YES	10
PRO	164	07-19	07-24	0.0115	ABR	NO	-	-	-
PRO	89	07-25	08-07	0.0110	INC	NO	-	-	-
PRO	76	12-19	12-24	0.0114	ABR	YES	374	NO	-
ROM	272	10-19	11-06	0.0114	ABR	NO	-	-	-
ROM	65	02-13	02-18	0.0115	ABR	YES	118	YES	9
ROM	165	04-25	04-30	0.0117	ABR	YES	118	YES	10
ROM	587	06-25	06-30	0.0118	ABR	YES	485	NO	-
ROM	57	07-25	07-30	0.0115	ABR	NO	-	-	-
ROM	356	07-31	08-07	0.0110	INC	YES	420	YES	10
ROM	239	10-01	10-16	0.0115	ABR	YES	485	YES	27
FOR	95	07-13	08-17	0.0010	INC	NO	-	-	-
FOR	97	11-13	12-06	0.0046	ABR	YES	175	YES	1

Table 6.2: Results of leak detection and localization across all districts of the Cavallermaggiore network. Columns “Det.” and “MAS” indicate whether the leak was detected and which sensor was identified as the Most Affected Sensor, respectively. Columns “Loc.” and “Rank” report whether the true leaking pipe appears in the ranked candidate list and its corresponding position within the list (with lower values indicating more accurate localization).

An additional consideration emerges when comparing abrupt (ABR) and incipient (INC) leaks. The algorithm does not exhibit a clear differentiation in performance between these two classes: both types are generally detected with similar reliability, and the localization accuracy remains consistent across them. This indicates that the hydraulic features extracted by the methodology are sufficiently robust to capture leak signatures regardless of whether the event manifests as a sudden burst or a gradually developing fault.

Finally, the analysis of the ranking outcomes shows that, for the correctly localized cases, the predicted leaking pipe typically appears within a reasonably narrow portion of the ranked list. The best possible outcome corresponds to Rank = 1, while the worst among the successful cases is Rank = 35. This range reflects the position of the true leaking pipe within the prioritized list generated by the algorithm, where candidate pipes are ordered according to their likelihood of being the leak location. Therefore, the rank list provides a practical measure of localization precision: lower ranks indicate highly accurate predictions, while higher ranks still denote meaningful hydraulic proximity even when not perfectly aligned with the ground truth.

Table 6.3 summarizes the performance of the leakage localization per district. It can be observed that the methodology achieves perfect localization on detected leaks in some districts (e.g., CON, MDP, FOR), whereas in other districts (e.g., ROM, PRO), a smaller fraction of all leaks are successfully localized. This highlights both the strength of the approach in correctly identifying detected leaks and the challenges remaining in achieving comprehensive coverage across all leak events, particularly in more complex network configurations or under sparse sensor deployments.

District	Localized/Detected	Localized/Total Leaks
CON	1.00	0.71
BRA	0.83	0.71
EUR	0.83	0.83
MDP	1.00	1.00
PRO	0.67	0.67
ROM	0.75	0.43
FOR	1.00	0.50

Table 6.3: Summary of leak localization per district. “Localized/Detected” shows the fraction of detected leaks that were correctly localized, while “Localized/ Total Leaks” shows the fraction of all leaks in the district that were correctly localized.

Chapter 7

Conclusion

This thesis has addressed the problem of leakage monitoring in urban water distribution networks by combining hydraulic modeling with data-driven and model-based detection techniques, with a focus on the municipal systems of Cavallermaggiore and Marene in the Province of Cuneo, Italy. The work has been developed in close collaboration with Alpi Acque and supported by detailed GIS Master data, SCADA measurements, and operational knowledge from the utility.

From a methodological perspective, Part I of the thesis has focused on the construction and calibration of hydraulic models using EPANET. Starting from GIS Master records, the network topology of Cavallermaggiore and Marene was reconstructed, including pipes, nodes, valves, tanks, wells, and pumping stations. Nodal demands were allocated by integrating yearly metered consumption data with time patterns derived from SCADA flow measurements at district inlets. Particular attention was devoted to tank geometry reconstruction, pump curve definition, and the implementation of realistic control logic based on tank levels. The resulting models were calibrated and validated against observed pressures, flows, and tank levels, and subsequently used to support the definition of District Metered Areas in Cavallermaggiore.

Part II has introduced two leakage monitoring strategies based on different data and modeling assumptions. The flow-based approach relies exclusively on SCADA flow sensor data and uses the DARTS framework to perform data-driven anomaly detection at the district level. While this method is effective for identifying abnormal conditions, it does not provide spatial localization of leakages within the network. For this reason, a second, model-driven pressure-based detection and localization approach was adopted. This approach builds on the calibrated EPANET model developed in Part I, generates synthetic leakage scenarios following the BattLeDIM framework, and applies the LILA algorithm for both pressure-based detection and spatial localization.

The work concludes with a set of results for both methods. The flow-based approach was evaluated using the historical data from the utility company. As

the historical data was not entirely reliable, and issues with the sensor data were identified, it was impossible to give a quantitative measure of the results. However, from a quantitative point of view, the method flow-based algorithm results were satisfactory. The pressure-based approach was not tested on real data, as pressure sensors with the required quantities are not present in the WDNs under study. It was then decided to evaluate the method following the same approach as BattleDIM, by simulating leakages in the Cavallermaggiore network using the hydraulic model. The results were satisfactory, proving that with enough and well placed pressure sensors, it is possible to reliably reduce the search area for leakages in the network.

This thesis demonstrates that the adoption of simulation and machine learning technologies by water utility companies has a positive impact in leakage monitoring, even when data sources are scarce, and that these technologies, together with proper instrumentation of the network can lead to a better management of our water resources.

Bibliography

- [1] ISTAT. *Censimento delle acque per uso civile*. Italian National Institute of Statistics. 2023. URL: <https://www.istat.it/it/archivio/273802>.
- [2] Stelios G. Vrachimis et al. «Battle of the Leakage Detection and Isolation Methods». In: *Journal of Water Resources Planning and Management* 148.12 (2022), p. 04022068. DOI: 10.1061/(ASCE)WR.1943-5452.0001601.
- [3] S. R. Mounce, J. B. Boxall, and J. Machell. «Development and verification of an online artificial intelligence system for detection of bursts and other abnormal flows». In: *Journal of Water Resources Planning and Management* 136.3 (2010), pp. 309–318. DOI: 10.1061/(ASCE)WR.1943-5452.0000030.
- [4] Riccardo Muradore et al. «Anomaly detection in water distribution networks using time series analysis». In: *Journal of Hydroinformatics* 25.2 (2023), pp. 312–329. DOI: 10.2166/hydro.2023.123.
- [5] M. Romano, Z. Kapelan, and D. A. Savić. «Automated detection of pipe bursts and other events in water distribution systems». In: *Journal of Water Resources Planning and Management* 140.4 (2014), pp. 457–467. DOI: 10.1061/(ASCE)WR.1943-5452.0000339.
- [6] U.S. Environmental Protection Agency. *Water Audits and Water Loss Control for Public Water Systems*. Tech. rep. EPA 816-F-13-002. EPA, 2021. URL: <https://www.epa.gov/waterutilityresponse/water-audits-and-water-loss-control>.
- [7] Allan Lambert and Wolfgang Hirner. «Losses from Water Supply Systems: Standard Terminology and Recommended Performance Measures». In: *IWA Blue Pages* (2000).
- [8] Stelios G. Vrachimis, Demetrios G. Eliades, and Marios M. Polycarpou. «CCWI 2020: The Battle of the Leakage Detection and Isolation Methods». In: *Proceedings of the 2nd International CCWI/WDSA Joint Conference*. Beijing, China, 2020.
- [9] Raido Puust et al. «A review of methods for leakage management in pipe networks». In: *Urban Water Journal* 7.1 (2010), pp. 25–45. DOI: 10.1080/15730621003610878.

- [10] Andrew F. Colombo, Philip Lee, and Bryan W. Karney. «A selective literature review of transient-based leak detection methods». In: *Journal of Hydro-environment Research* 2.4 (2009), pp. 212–227. DOI: 10.1016/j.jher.2009.02.003.
- [11] Roland Liemberger and Alan Wyatt. «Quantifying the global non-revenue water problem». In: *Water Science and Technology: Water Supply* 19.1 (2007), pp. 1–14. DOI: 10.2166/ws.2018.129.
- [12] Allan Lambert. «What do we know about pressure: leakage relationships in distribution systems?» In: *IWA Conference on System Approach to Leakage Control and Water Distribution Systems Management*. Brno, Czech Republic, 1999.
- [13] Jamal Maher Alkassseh et al. «Applying Minimum Night Flow to Estimate Water Loss Using Statistical Modeling: A Case Study in Kinta Valley, Malaysia». In: *Water Resources Management* 27.5 (2013), pp. 1439–1455. DOI: 10.1007/s11269-012-0247-2.
- [14] Malcolm Farley and Stuart Trow. «Losses in Water Distribution Networks: A Practitioners’ Guide to Assessment, Monitoring and Control». In: *IWA Publishing* (2008).
- [15] Lewis A. Rossman. *EPANET 2.2 User Manual*. EPA/600/R-20/133. U.S. Environmental Protection Agency. 2020. URL: <https://epanet2.readthedocs.io/>.
- [16] Lewis A. Rossman. «EPANET 2: Users Manual». In: *Water Supply and Water Resources Division, National Risk Management Research Laboratory* (2000).
- [17] Z. Y. Wu, P. Sage, and D. Turtle. «Pressure-dependent leak detection model and its application to a district water system». In: *Journal of Water Resources Planning and Management* 136.1 (2010), pp. 116–128. DOI: 10.1061/(ASCE)0733-9496(2010)136:1(116).
- [18] Orazio Giustolisi, Dragan Savic, and Zoran Kapelan. «Pressure-driven demand and leakage simulation for water distribution networks». In: *Journal of Hydraulic Engineering* 134.5 (2008), pp. 626–635. DOI: 10.1061/(ASCE)0733-9429(2008)134:5(626).
- [19] Janet M. Wagner, Uri Shamir, and David H. Marks. «Water distribution reliability: Simulation methods». In: *Journal of Water Resources Planning and Management* 114.3 (1988), pp. 276–294. DOI: 10.1061/(ASCE)0733-9496(1988)114:3(276).
- [20] Dragan A. Savic, Zoran Kapelan, and Pauline M. R. Jonkergouw. «Quo vadis water distribution model calibration?» In: *Urban Water Journal* 6.1 (2009), pp. 3–22. DOI: 10.1080/15730620802613380.

- [21] Zoran S. Kapelan, Dragan A. Savic, and Godfrey A. Walters. «Multiobjective design of water distribution systems under uncertainty». In: *Water Resources Research* 41.11 (2005), W11407. DOI: 10.1029/2004WR003787.
- [22] Kalyanmoy Deb et al. «A fast and elitist multiobjective genetic algorithm: NSGA-II». In: *IEEE Transactions on Evolutionary Computation* 6.2 (2002), pp. 182–197. DOI: 10.1109/4235.996017.
- [23] Bryan A. Tolson and Christine A. Shoemaker. «Dynamically dimensioned search algorithm for computationally efficient watershed model calibration». In: *Water Resources Research* 43.1 (2007), W01413. DOI: 10.1029/2005WR004723.
- [24] Christopher J. Hutton et al. «Dealing with uncertainty in water distribution system models: A framework for real-time modeling and data assimilation». In: *Journal of Water Resources Planning and Management* 140.2 (2014), pp. 169–183. DOI: 10.1061/(ASCE)WR.1943-5452.0000325.
- [25] Marco Salvetti, Matteo Giuliani, and Andrea Castelletti. «A coupled modelling framework for the optimal design of multi-purpose water reservoir systems». In: *Journal of Hydroinformatics* 21.5 (2019), pp. 851–865. DOI: 10.2166/hydro.2019.140.
- [26] Doosun Kang and Kevin Lansey. «Real-time demand estimation and confidence limit analysis for water distribution systems». In: *Journal of Hydraulic Engineering* 135.10 (2009), pp. 825–837. DOI: 10.1061/(ASCE)HY.1943-7900.0000086.
- [27] Jamie Morrison, Stuart Tooms, and Duncan Rogers. «District metered areas: Guidance notes». In: *Water Loss Task Force, IWA* (2007).
- [28] Dominic L. Boccelli et al. «Optimal scheduling of booster disinfection in water distribution systems». In: *Journal of Water Resources Planning and Management* 124.2 (2017), pp. 99–111. DOI: 10.1061/(ASCE)0733-9496(1998)124:2(99).
- [29] Carlo Giudicianni et al. «Topological Taxonomy of Water Distribution Networks». In: *Water* 10.4 (2018), p. 444. DOI: 10.3390/w10040444.
- [30] S. Russell, V. Babovic, and S. Dogruel. «A unified framework for real-time optimal control of water distribution networks». In: *Water Supply* 19.3 (2019), pp. 862–871. DOI: 10.2166/ws.2018.143.
- [31] Dalius Misiunas et al. «Pipeline break detection using pressure transient monitoring». In: *Journal of Water Resources Planning and Management* 131.4 (2005), pp. 316–325. DOI: 10.1061/(ASCE)0733-9496(2005)131:4(316).
- [32] E. S. Page. «Continuous inspection schemes». In: *Biometrika* 41.1/2 (1954), pp. 100–115. DOI: 10.2307/2333009.

- [33] Michèle Basseville and Igor V. Nikiforov. *Detection of Abrupt Changes: Theory and Application*. Prentice Hall, 1993. ISBN: 978-0131267800.
- [34] Sean J. Taylor and Benjamin Letham. «Forecasting at scale». In: *The American Statistician* 72.1 (2018). Prophet forecasting algorithm, pp. 37–45. DOI: 10.1080/00031305.2017.1380080.
- [35] Julien Herzen et al. «Darts: User-Friendly Modern Machine Learning for Time Series». In: *Journal of Machine Learning Research*. Vol. 23. 124. 2022, pp. 1–6. URL: <https://jmlr.org/papers/v23/21-1177.html>.
- [36] Sepp Hochreiter and Jürgen Schmidhuber. «Long short-term memory». In: *Neural Computation* 9.8 (1997), pp. 1735–1780. DOI: 10.1162/neco.1997.9.8.1735.
- [37] Antonio Candelieri and Francesco Archetti. «Identifying typical urban water demand patterns for a reliable short-term forecasting—The icewater project approach». In: *Procedia Engineering* 119 (2015), pp. 1004–1012. DOI: 10.1016/j.proeng.2015.08.930.
- [38] Marcos V. Casillas, Luis E. Garza-Castañón, and Vicenç Puig Cayuela. «Model-based leak detection and location in water distribution networks considering an extended-horizon analysis of pressure sensitivities». In: *Journal of Hydroinformatics* 15.3 (2013), pp. 834–851. DOI: 10.2166/hydro.2013.019.
- [39] Gustavo Sanz and Ramón Pérez. «Leak detection and localization through demand components calibration». In: *Journal of Water Resources Planning and Management* 142.2 (2016), p. 04015057. DOI: 10.1061/(ASCE)WR.1943-5452.0000592.
- [40] Adrià Soldevila et al. «Leak localization in water distribution networks using a mixed model-based/data-driven approach». In: *Control Engineering Practice* 55 (2016), pp. 162–173. DOI: 10.1016/j.conengprac.2016.07.006.
- [41] Stelios Sophocleous, Dragan Savić, and Zoran Kapelan. «Leak localization in a real water distribution network based on search-space reduction». In: *Journal of Water Resources Planning and Management* 145.7 (2019), p. 04019024. DOI: 10.1061/(ASCE)WR.1943-5452.0001079.
- [42] Filippo Pecci, Ivan Stoianov, and Avi Ostfeld. «Combining differential evolution and gradient-based algorithms for multi-objective optimization of water distribution systems». In: *Journal of Water Resources Planning and Management* 145.5 (2019), p. 04019010. DOI: 10.1061/(ASCE)WR.1943-5452.0001051.
- [43] Willem B. C. de Schaetzen, Dragan A. Savic, and Godfrey A. Walters. «Optimal sampling design for model calibration using shortest path, genetic and entropy algorithms». In: *Urban Water* 2.2 (2000), pp. 141–152. DOI: 10.1016/S1462-0758(00)00052-5.

- [44] John Mashford et al. «An approach to leak detection in pipe networks using analysis of monitored pressure values by support vector machine». In: *Proceedings of the International Conference on Network and System Security* (2009), pp. 534–539. DOI: 10.1109/NSS.2009.101.
- [45] Ivo Daniel et al. «A Sequential Pressure-Based Algorithm for Data-Driven Leakage Identification and Model-Based Localization in Water Distribution Networks». In: *Journal of Water Resources Planning and Management* 148.6 (June 2022). ISSN: 1943-5452. DOI: 10.1061/(asce)wr.1943-5452.0001535. URL: [http://dx.doi.org/10.1061/\(ASCE\)WR.1943-5452.0001535](http://dx.doi.org/10.1061/(ASCE)WR.1943-5452.0001535).
- [46] Ramon Sarrate, Fatiha Nejari, and Albert Rosich. «Sensor placement for fault diagnosis performance maximization in distribution networks». In: *Control Engineering Practice* 31 (2014), pp. 1–10. DOI: 10.1016/j.conengprac.2014.06.002.
- [47] David B. Steffebauer and Daniela Fuchs-Hanusch. «Efficient sensor placement for leak localization considering uncertainties». In: *Water Resources Management* 30.14 (2016), pp. 5517–5533. DOI: 10.1007/s11269-016-1504-6.
- [48] Marcos V. Casillas et al. «Optimal sensor placement for leak location in water distribution networks using genetic algorithms». In: *Sensors* 13.11 (2013), pp. 14984–15005. DOI: 10.3390/s131114984.
- [49] Jessica Chicco et al. «Dual-use value of network partitioning for water system management and protection from malicious contamination». In: *Journal of Hydroinformatics* 21.4 (2019), pp. 619–637. DOI: 10.2166/hydro.2019.078.
- [50] Alpiacque S.p.A. *GIS Master Technical Design Platform*. Geographic Information System for water network infrastructure management. 2024. URL: <https://cloud-farm.gismaster.it/>.