

Leak Localization in Water Distribution Networks Using Data-Driven and Model-Based Approaches

Luis Romero-Ben¹, Débora Alves², Joaquim Blesa³, Gabriela Cembrano⁴, Vicenç Puig⁵, and Eric Duviella⁶

¹Institut de Robòtica i Informàtica Industrial (CSIC-UPC). Carrer Llorens Artigas, 4-6, 08028 Barcelona, Spain. Email: luis.romero.ben@upc.edu (corresponding author)

²Supervision, Safety and Automatic Control Research Center (CS2AC) of the Universitat Politècnica de Catalunya, Campus de Terrassa, Gaia Building, Rambla Sant Nebridi, 22, 08222 Terrassa, Barcelona, Spain. Email: adeboracris@gmail.com

³Dr.Eng. Institut de Robòtica i Informàtica Industrial (CSIC-UPC). Carrer Llorens Artigas, 4-6, 08028 Barcelona, Spain; Serra Húnter fellow, Automatic Control Department of the Universitat Politècnica de Catalunya, Avinguda Eduard Maristany, 16, 08019, Barcelona, Spain. Email: joaquim.blesa@upc.edu

⁴Dr.Eng. Institut de Robòtica i Informàtica Industrial (CSIC-UPC). Carrer Llorens Artigas, 4-6, 08028 Barcelona, Spain. Email: gabriela.cembrano@upc.edu

⁵Dr.Eng. Institut de Robòtica i Informàtica Industrial (CSIC-UPC). Carrer Llorens Artigas, 4-6, 08028 Barcelona, Spain; Supervision, Safety and Automatic Control Research Center (CS2AC) of the Universitat Politècnica de Catalunya, Campus de Terrassa, Gaia Building, Rambla Sant Nebridi, 22, 08222 Terrassa, Barcelona, Spain. Email: vicenc.puig@upc.edu

⁶Dr.Eng. Informatics and Automatics Department, IMT Lille Douai, Lille, France. Email: eric.duviella@imt-lille-douai.fr

ABSTRACT

The detection and localization of leaks in water distribution networks (WDNs) is one of the major concerns of water utilities, due to the necessity of an efficient operation that satisfies the

25 worldwide growing demand of water. There exists a wide range of methods, from equipment-based
26 techniques that rely only on hardware devices to software-based methods that exploit models and
27 algorithms as well. Model-based approaches provide an effective performance but they rely on
28 the availability of an hydraulic model of the WDN, while data-driven techniques only require
29 measurements from the network operation although they may produce less accurate results. This
30 paper proposes two methodologies: a model-based approach that uses the hydraulic model of the
31 network, as well as pressure and demand information; and a fully data-driven method based on
32 graph interpolation and a new candidate selection criteria. Their complementary application was
33 successfully applied to the Battle of the Leakage Detection and Isolation Methods (BattLeDIM)
34 2020 challenge, and the achieved results are presented in this paper to demonstrate the suitability
35 of the methods.

36 INTRODUCTION

37 A recent study, developed in [Liemberger and Wyatt \(2019\)](#), estimated that leaks account for up
38 to 126 billion cubic meters of water per year worldwide (expressed as non-revenue water), which
39 represents a remarkably significant quantity considering the worldwide growing demand, supposed
40 to increase by 55% between 2000 and 2050 according to [Leflaive \(2012\)](#). Furthermore, apart from
41 the associated economical and operational costs, water leaks increase the risk of contamination
42 ([Xu et al. 2014](#)) and health problems ([LeChevallier et al. 2003](#)). Multitude of solutions have been
43 proposed during the years to address the leak detection and localization problem (see [Chan et al.](#)
44 [\(2018\)](#) for an extensive review). They are typically classified into two categories: hardware-based
45 and software-based methods ([Li et al. 2015](#)).

46 On the one hand, hardware-based methods use hardware devices to detect the existence and
47 position of bursts. They are usually divided into acoustic methods, including listening rods, leak
48 correlators and leak noise loggers ([Mutikanga et al. 2013](#)); and non-acoustic approaches, like gas
49 injection, ground penetrating radar technology or thermal infrared imaging among others ([Fanner](#)
50 [et al. 2007](#)). Even if these methods provide a high degree of accuracy, their usage is usually
51 prohibitive for large pipe networks due to its high associated costs and reduced detection range,

52 and hence their application is limited to small zones of the WDNs (Rajeswaran et al. 2018).

53 On the other hand, software-based techniques rely on models or algorithms that exploit ad-
54 ditional information from metering devices (pressure meters, flow sensors, etc.) to perform the
55 detection/localization of the leaks. These methods can be split into three main categories: transient-
56 based, model-based and data-driven methods.

57 Transient-based approaches analyse the transients induced by leaks using signal processing
58 techniques. A leak detection method that exploits transient-based information and genetic algo-
59 rithms is described in Vítkovský et al. (2000). The detection task is tackled in Kapelan et al.
60 (2003) by a hybrid inverse transient procedure, formulated as a constrained optimization problem
61 of a weighted least-squares cost function. Besides, Wang et al. (2020) proposes a method that
62 uses matched-field processing to locate leaks by incorporating prior information of the modelling
63 error. Techniques based on leak transients are also used for leak diagnosis in pipelines, as shown
64 in Pérez-Pérez et al. (2021) and Torres et al. (2021).

65 Model-based methodologies use hydraulic models and simulators, calibrating both the network
66 characteristics and the demands, to compare simulated hydraulic information with actual measure-
67 ments from the WDN (see the introduction in Sanz et al. (2016) for a review). A leak detection
68 and localization technique using flow velocities is presented in Goulet et al. (2013), built as an
69 error-domain model falsification methodology. A widely-used localization approach that matches
70 pressure residuals to a fault signature matrix obtained by means of hydraulic simulations is proposed
71 in Pérez et al. (2014). A similar consideration is used in Sophocleous et al. (2019), performing a
72 sensitivity analysis along with a search-space reduction approach to find the leakage location.

73 While model-based approaches have been widely researched and exploited due to their efficiency
74 and effectiveness (Duan et al. 2011), the associated performance is limited by the difficulty in the
75 selection and calibration of the corresponding mathematical models (Menapace et al. 2018), the
76 diversity and complexity of WDNs (Kim et al. 2016), and the presence of modelling errors like nodal
77 demand uncertainties and measurement noise (Blesa and Pérez 2018). Most of these disadvantages
78 may be gradually overcome using data-driven and machine learning techniques (Ferrandez-Gamot

79 [et al. 2015](#)) and their reduced or non-existent dependency on an hydraulic model.

80 Data-driven strategies analyse the measurements from monitoring devices, mining knowledge to
81 detect leaks and identify their location (see [Wu and Liu \(2017\)](#) for an extensive review). Information
82 from accelerometers is exploited in [Kang et al. \(2017\)](#) to feed a two-phase method supported by a
83 convolutional neural network (CNN) - support vector machine (SVM) architecture that detect leaks
84 and a graph method based on virtual nodes that locates them. The concept of Kantorovich distance
85 is applied in [Arifin et al. \(2018\)](#) to detect and locate leaks by exploiting the pipeline leak signature
86 and identifying possible changes in the pipeline status using mass flow rates and pressure data. An
87 efficient multistage method is presented in [Huang et al. \(2020\)](#), which uses valve operations (VOs)
88 to split the demand metering area (DMA) into two zones and identify the leak location within these
89 regions by means of a water balance analysis based on smart demand meters.

90 Recently proposed data-driven approaches use pressure sensors due to their the lower cost and
91 easier installation in comparison with other kind of meters, resulting in an attractive option for water
92 utilities ([Soldevila et al. 2021](#)). In [Han et al. \(2018\)](#), a two-stage strategy is used to estimate the
93 WDN state by means of a Gauss-Newton Belief Propagation inference scheme applied to hydraulic
94 heads at certain nodes of the network, and a clustering method to decompose the WDN and isolate
95 the leak. A deep-learning scheme is proposed in [Zhou et al. \(2019\)](#) to locate leaks using pressure
96 meters that are placed at limited, optimised places for a short period. A data-driven approach is
97 developed in [Soldevila et al. \(2021\)](#), interpolating the pressure at every node of the network from
98 certain measured values by means of the Kriging interpolation technique, comparing leak and
99 leak-free scenarios to locate the leak and using Dempster-Shafer reasoning to deal with uncertainty.

100 This article proposes two different and complementary techniques to locate leaks: a model-
101 based method that uses a hydraulic model of the WDN and pressure and demand measurements
102 and/or well-calibrated demand estimations; and a data-driven approach that only requires pressure
103 information from some inner nodes and the topology of the network. Both of them are applied to
104 the Battle of the Leakage Detection and Isolation Methods (BattLeDIM) 2020 challenge ([Vrachimis
105 et al. 2020](#)), using them in a complementary manner to improve the performance.

106 Both methodologies present several contributions with respect to previous approaches:

- 107 • They handle the multi-leak problem, overcoming the classical hypothesis of the appearance
108 of a single leak at a time (Goulet et al. 2013; Pérez et al. 2014), which is assumed in most
109 state-of-the-art techniques (Soldevila et al. 2016).
- 110 • Regarding the data-driven method, it reduces the complexity of the interpolation stage
111 by using a quadratic programming approach that exploits the topology of the network,
112 concretely the length of the pipes. Other interpolation-based methods like Soldevila et al.
113 (2021) require extra information to be known (like diameters) or estimated during the
114 interpolation (distribution of flow in pipes, pipe roughness, etc.). Moreover, in the referred
115 work, a basic residual computation is employed to select the candidate to be the origin
116 of the leak. This approach is completed in the method proposed here by combining the
117 information of both the basic residuals and the relation among the hydraulic heads of all the
118 nodes of the network (interpolated and measured) in a single metric.

119 Furthermore, a leak detection and estimation technique is proposed to complete the solution of
120 the challenge and feed the localization strategies with the leak appearance time instants. Then, the
121 proposed localization methods depend on the proper operation of the detection stage.

122 **METHODOLOGY**

123 Water distribution networks across the world present different characteristics in structure, size,
124 demand patterns, components, etc. Moreover, the distribution, amount and properties of the
125 measuring devices installed throughout the networks varies from one site to another. This fact
126 indicates the necessity of adapting the selection and usage of software-based techniques.

127 Thus, two complementary approaches are proposed in the following to locate leaks in WDNs:

- 128 • A model-based methodology that exploits the existence of a well-calibrated hydraulic model
129 of the WDN, as well as the availability of reliable pressure and demand information.
- 130 • A fully data-driven technique that only requires minimal topological knowledge of the
131 network and measurements from pressure sensors distributed at a set of inner nodes.

132 The model-based method faces difficulties when there is a lack of demand measurements or
 133 estimations, while the data-driven strategy requires a minimum pressure sensor density in order
 134 to operate, as it is its only source of hydraulic information. The WDNs can be divided into areas
 135 with different features (like in the presented case study), e.g., sensorization properties (amount,
 136 placement, precision and type of installed sensors), existence of weirs, tanks, valves... Therefore, the
 137 best option between these two methods can be selected depending on the availability of information
 138 and/or model of the WDN, and they can be even used in a complementary manner, so that the
 139 weaknesses of one method are compensated by the strengths of the other.

140 **Leak detection and estimation**

141 The proposed method uses sensor fusion calculations to analyse the flow of water supplied to
 142 the DMA at all hours of the day, not only during the night hours. Initially, it is assumed that the
 143 demand forecasting method is calibrated with historical data from the DMAs (Donkor et al. 2014),
 144 which will provide a good approximation of the current inflow y at time k :

$$145 \quad y(k) = \hat{y}(k) + e(k) \quad (1)$$

146 where $k = 0, 1, 2, \dots$, denotes the discrete time corresponding to time $0, T_s, 2T_s, \dots$ (T_s is the sample
 147 time of the demand forecasting model), $\hat{y}(k)$ is the demand forecast and $e(k)$ is the error that, for
 148 this study, is considered to be adjusted by a normal distribution (Malik 2016). As the incoming
 149 demand is more accurate in some periods of the day, a periodic variation in time T is considered:

$$150 \quad e(k) \sim N(0, \sigma^2(k)) \quad \text{with} \quad \sigma^2(k) = \sigma^2(k + T) \quad (2)$$

151 In the case of the presence of a leak, i.e., $f(k) > 0$, Equation (1) leads to:

$$152 \quad y(k) = \hat{y}(k) + e(k) + f(k) \quad \longrightarrow \quad \hat{f}(k) = y(k) - \hat{y}(k) = f(k) + e(k) \quad (3)$$

153 where $\hat{f}(k)$ is an approximation of the leak size given by the difference between the actual and the
 154 estimated inlet flow, with a leak estimation error equal to the demand forecasting error.

155 Given the current inlet flow value and previous values in a time window with N_W samples,
 156 Equation (3) can be applied considering N_W different leak estimations, approximating the leak as:

$$157 \quad f(k) \approx \bar{f}(k) = \sum_{i=0}^{N_W-1} \frac{f(k-i)}{N_W} \quad (4)$$

158 An average leak estimation $\hat{f}(k)$ can be computed at instant k applying the maximum likelihood
 159 estimation method to the joint probability distribution of the N_W estimations fused in $\bar{f}(k)$:

$$160 \quad \hat{f}(k) = \frac{\sum_{i=0}^{N_W-1} \frac{\hat{f}(k-i)}{\sigma^2(k-i)}}{\sum_{i=0}^{N_W-1} \frac{1}{\sigma^2(k-i)}} \quad (5)$$

161 In a non-leak scenario, $\hat{f}(k)$ will lead to small values (but different from zero) due to demand
 162 estimation errors, whereas its value will increase in a leak scenario. Thus, the leak detection is
 163 formulated as a change detection problem that can be solved by computing a threshold λ that will
 164 determine the value of $\hat{f}(k)$ above which it can be assumed that a leak exists. This value can be
 165 computed applying Equation (5) to historical leak-free data, considering the worst-case scenario λ
 166 to be equal to the maximum value of $\hat{f}(k)$ computed for the whole historical non-leak data.

167 Then, the leak detection is triggered as stated by the following expression:

$$168 \quad \begin{cases} \hat{f}(k) > \lambda \Rightarrow \text{Leak} \\ \text{Otherwise} \Rightarrow \text{No Leak} \end{cases} \quad (6)$$

169 **Data-driven methodology for leak localization**

170 The proposed fully data-driven approach is based on two phases: a graph-based interpolation
 171 that approximates the complete hydraulic state of the network and a geometric comparison between
 172 the states in leak and leak-free scenarios. This approach removes the need for a hydraulic model,
 173 as it only uses measured pressure values and the WDN topology. It presents the following features:

- 174 • It is applicable to the measurements of a single time-instant, but it can be easily extended
 175 to integrate temporal information from time-series.

- 176 • The leak-free data can be obtained from a historical dataset of hydraulic measurements,
177 provided by the water utility. In order to deal with possible differences in the demand
178 conditions between the leak and leak-free scenarios, updated nominal information can be
179 obtained from the previous days to the appearance of the leak. Moreover, these demand
180 discrepancies can be reduced even more by selecting nominal scenarios with an analogous
181 consumption pattern, e.g., the same day of the week. In this way, an hydraulic model is not
182 required to obtain the nominal data.
- 183 • It provides a set of leak candidate inner nodes. The size of the set depends on the sensors
184 limitations and network structure. However, it offers the possibility of approximating the
185 exact leak location with sufficient reliability.
- 186 • It is flexible regarding the network structure, so that modifications of the WDN only imply
187 the update of the topological information. Besides, it does not require demand information,
188 which is one of the critical points of most of the model-based strategies.

189 This leak localization method must be used together with the presented detection technique.

190 *Graph-based state interpolation*

191 In this work, the location of leaks is determined through the comparison between the actual
192 hydraulic state of the network and a leak-free reference. The hydraulic heads at the nodes are
193 chosen as a representative state variables due to the advantages of pressure sensors and the usage
194 of this information by water utilities to determine the availability of water service.

195 In order to minimize the installation of sensors, the complete state is interpolated from a reduced
196 set of measurements. The non-linear relation among the hydraulic heads of neighbouring nodes,
197 typically described by the Hazen-Williams equation, is approximated considering the state of a
198 node to be computed as the weighted linear contribution of its neighbours, hence simplifying the
199 interpolation procedure and expanding the set of tools for the derivation of the network state.

200 Let us model the network structure by means of a simple directed graph referred to as
201 $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where the node set \mathcal{V} is representing the set of junctions of the WDN and the

202 edge set \mathcal{E} stands for the set of pipes. The total number of junctions of the network, i.e. the
 203 cardinality $|\mathcal{V}|$ of the node set, is referred to as n , as well as the number of pipes, i.e., the car-
 204 dinality $|\mathcal{E}|$ of the edge set, is represented as m . A node is referred to as $v_i \in \mathcal{V}$, and an edge
 205 $e_{ij} = (v_i, v_j) \in \mathcal{E}$ connects source node v_i with sink node v_j , so that they are its endpoints. A set
 206 of important matrices can be extracted from the graph characteristics:

- 207 • The connectivity among nodes is used to derive the adjacency matrix $A(\mathcal{G}) \in \mathbb{R}^{n \times n}$:

$$208 \quad a_{ij} = \begin{cases} 1, & \text{if } e_{ij} \in \mathcal{E} \text{ or } e_{ji} \in \mathcal{E} \\ 0, & \text{otherwise} \end{cases} \quad (7)$$

- 209 • The directionality of the graph edges represents the direction of the water flows through the
 210 network pipes. Considering that v_i is the source of an edge e_o (or e_{ij}), while v_j is its sink,
 211 the directionality is represented by the incidence matrix $B(\mathcal{G}) \in \mathbb{R}^{m \times n}$ as follows:

$$212 \quad b_{oj} = \begin{cases} 1, & \text{if } e_o = (v_i, v_j) \in \mathcal{E} \\ -1, & \text{if } e_o = (v_j, v_i) \in \mathcal{E} \\ 0, & \text{otherwise} \end{cases} \quad (8)$$

213 with o from 1 to m indexing the edges of the graph. An approximated $B(\mathcal{G})$ can be estimated
 214 considering the source-sink relation among the reservoirs and all the inner nodes, as the
 215 exact flow direction of each pipe may be unknown and/or it may vary.

- 216 • An edge e_{ij} is associated to a cost value, related in this case to the pipe length l_{ij} . The edge
 217 costs are exploited to generate a weighted adjacency matrix $W(\mathcal{G}) \in \mathbb{R}^{n \times n}$:

$$218 \quad w_{ij} = \begin{cases} \frac{1}{l_{ij}}, & \text{if } l_{ij} \neq 0 \\ 0, & \text{otherwise} \end{cases} \quad (9)$$

219 where the definition of the weight w_{ij} is designed so that closer neighbours affect the value
 220 of the considered node in a higher degree.

- 221 • The degree matrix $D(\mathcal{G}) \in \mathbb{R}^{n \times n}$ is derived from $W(\mathcal{G})$, i.e., $d_{ij} = \sum_{h=1}^n w_{ih}$ only if $i = j$;
 222 being zero otherwise (h plays the role of the columns index in w_{ih} because j is used in d_{ij}).

223 With the aid of the these matrices, the state x_i of a certain node v_i can be expressed as:

$$224 \quad x_i = \frac{1}{d_i} \mathbf{w}_i \mathbf{x} \quad (10)$$

225 where $d_i = d_{ii} = d_{ij}$, $\mathbf{w}_i \in \mathbb{R}^{1 \times n}$ denotes the i -th row of $W(\mathcal{G})$, and $\mathbf{x} \in \mathbb{R}^n$ represents the complete
 226 state vector. The estimation of \mathbf{x} can be achieved through the minimization of the sum of the
 227 quadratic difference between each node actual value and the one computed by Equation (10):

$$228 \quad \sum_{i=1}^n \left[x_i - \frac{1}{d_i} \mathbf{w}_i \mathbf{x} \right]^2 = (\mathbf{x} - D^{-1} W \mathbf{x})^T (\mathbf{x} - D^{-1} W \mathbf{x}) = \mathbf{x}^T (I_n - D^{-1} W)^T (I_n - D^{-1} W) \mathbf{x} =$$

$$\mathbf{x}^T (D^{-1} (D - W))^T (D^{-1} (D - W)) \mathbf{x} = \mathbf{x}^T (D^{-1} L)^T (D^{-1} L) \mathbf{x} = \mathbf{x}^T L D^{-2} L \mathbf{x} \quad (11)$$

229 where $D = D(\mathcal{G})$ and $W = W(\mathcal{G})$, and $L = L^T = L(\mathcal{G}) = D(\mathcal{G}) - W(\mathcal{G})$ is the unnormalized
 230 Laplacian matrix of graph \mathcal{G} (Mohar et al. 1991). I_n stands for the identity matrix of size n .

231 In order to provide the available hydraulic information to the minimization process, the state of
 232 the metered nodes is filled using the measured hydraulic heads through an equality constraint:

$$233 \quad S \mathbf{x} = \mathbf{x}^s \quad (12)$$

234 where $S \in \mathbb{R}^{n \times n}$ is a diagonal matrix which has a value of 1 at its $i - i$ component if there exists a
 235 sensor at v_i , and a value of 0 otherwise. Vector $\mathbf{x}^s \in \mathbb{R}^n$ contains the head values of the measured
 236 nodes at the corresponding components, while the rest of values are 0. The pressure at the water
 237 inputs of the network is supposed to be known, as it is commonly available in most of the WDNs.

238 The approximation in Equation (10) denotes the harmonic property of functions in graphs,
 239 explained in Zhu et al. (2003) to interpolate the value of a function to unknown vertices following
 240 a smoothing strategy. In that work, the harmonic property is implicitly pursued by minimizing a

241 weighted quadratic energy function over the function values. However, the harmonic property is
 242 explicitly pursued in Equation (11) by directly minimizing the quadratic difference between the
 243 actual node state and the value estimated by computing the average of the state in neighbouring
 244 nodes. This cost function does not impose a smoothing objective like the one in [Zhu et al. \(2003\)](#),
 245 so all the possible combinations of neighbouring states that produce the same state for a certain
 246 node are equally considered.

247 This provides a higher degree of freedom in the solution to attain a second objective, related to
 248 the directionality of the flows. In WDNs, the direction of water flow through a pipe is determined
 249 by the sign of the difference in hydraulic head between its corresponding junctions, taking into
 250 account that water flows in the direction of decreasing hydraulic head. This fact can be translated
 251 to the relation among the states of the nodes of graph \mathcal{G} by means of the following inequality:

$$252 \quad B\mathbf{x} \leq \boldsymbol{\gamma} \quad (13)$$

253 where $B = B(\mathcal{G})$ is defined as stated in Equation (8) and $\boldsymbol{\gamma} \in \mathbb{R}^n$ stands for a vector with a value of
 254 γ for all its components. This expression would imply a requirement about the difference of state
 255 (that is, approximated hydraulic head) between adjacent nodes, so that it must be lower or equal to
 256 a certain threshold γ . This value is introduced because $B(\mathcal{G})$ is estimated from the WDN topology
 257 and hence some slack is desirable in the fulfilment of the directionality constraints. However, this
 258 threshold is included in a second objective of the optimization cost function to compute its lowest
 259 positive feasible value. Thus, the complete optimization problem can be arranged as:

$$\begin{aligned}
 & \min_{\mathbf{x}} \quad \frac{1}{2} [\mathbf{x}^T L D^{-2} L \mathbf{x} + \alpha \gamma^2] \\
 & \text{s.t.} \quad B\mathbf{x} \leq \boldsymbol{\gamma} \\
 & \quad \quad \gamma > 0 \\
 & \quad \quad S\mathbf{x} = \mathbf{x}^s
 \end{aligned}
 \quad (14)$$

261 where α allows to settle the importance of the directionality objective.

262 *Leak candidate selection method*

263 As leak and leak-free scenarios must be compared to locate the leak, the graph-based state
 264 interpolation is applied to both the actual and nominal hydraulic information, i.e., the measured
 265 hydraulic heads at a certain time instant of the detected leaky event and the ones at a leak-free time
 266 instant (with similar boundary conditions (Pérez et al. 2011)). The former are stored in a vector
 267 denoted as $\mathbf{x}_{leak} \in \mathbb{R}^n$, while the latter are collected to form a vector referred to as $\mathbf{x}_{nom} \in \mathbb{R}^n$.

268 The proposed leak candidate selection method considers the components of those vectors to
 269 represent the coordinates of n points in \mathbb{R}^2 (\mathbf{x}_{nom} provides the x-coordinates and \mathbf{x}_{leak} provides
 270 the y-coordinates). In the case of comparing two healthy scenarios, these points should be rather
 271 aligned because they would be only affected by the boundary conditions. However, the presence of
 272 a leak would alter the \mathbb{R}^2 location of the points representing the affected nodes, considering that a
 273 leak produces a reduction of the expected hydraulic head (Adedeji et al. 2017).

274 Thus, the objective of this stage consists of providing the set of most distant points from
 275 their predicted position, given by the best fitting line of the complete set of points. As this line
 276 is computed considering all of them, the distance between a certain node and the line not only
 277 encodes information about the change of state in that node, but also the relation among the state of
 278 all the nodes of the network. The best-fitting line can be expressed as $\mathbf{x}_{leak} = X_{nom}\boldsymbol{\phi}$, with $X_{nom} =$
 279 $[\mathbf{x}_{nom} \mathbf{1}_n] \in \mathbb{R}^{n \times 2}$ ($\mathbf{1}_n$ is a column vector whose components are all 1) and $\boldsymbol{\phi} = [\phi_1 \phi_2]^T \in \mathbb{R}^2$. The
 280 latter is the vector containing the line parameters, which are computed by solving the least-squares
 281 problem, i.e., $\boldsymbol{\phi} = (X_{nom}^T X_{nom})^{-1} X_{nom}^T \mathbf{x}_{leak}$.

282 The perpendicular distance from the set of points to the best-fitting line can be computed as:

$$283 \quad \boldsymbol{\delta} = [\mathbf{x}_{nom} \ \mathbf{x}_{leak} \ \mathbf{1}_n] \left[\begin{array}{ccc} \frac{\phi_1}{\sqrt{\phi_1^2 + 1}} & \frac{-1}{\sqrt{\phi_1^2 + 1}} & \frac{\phi_2}{\sqrt{\phi_1^2 + 1}} \end{array} \right]^T \quad (15)$$

284 so $\boldsymbol{\delta} \in \mathbb{R}^n$ stands for the vector of computed distances. The sign is kept in the calculation since
 285 only points with a positive distance value can be leak candidates, i.e., points that are located below
 286 the best-fitting line because their y-coordinate value, that is the leaky one, is lower than expected.

287 A criterion must be selected to pick a set of nodes from \mathcal{V} depending on the values at δ . In this
288 case, the standard deviation σ_δ of the distance vector is computed to play the role of a threshold,
289 i.e., a node v_i must have an associated distance value δ_i that exceeds σ_δ in order to consider this
290 node as a candidate to be the leak origin. The final candidates can be ordered from most to least
291 probable by means of their associated distance δ_i , and therefore the node that corresponds to the
292 highest distance value is considered as the best candidate.

293 **Model-based methodology for leak localization**

294 The proposed model-based leak localization method uses a hydraulic simulator to simulate
295 theoretical pressure values caused by all potential leaks once a leak has been detected and its
296 magnitude has been estimated. Simulated pressure values at different leak locations (hypothesis)
297 are compared with the DMA measured pressure values to determine the most probable leak location.
298 After the leak localization procedure, the hydraulic simulator is updated with the new extra demand,
299 whose magnitude is the leak estimation value at the leak localization. Thus, the proposed method
300 can tackle the problem of multi-leaks $f_{j_1}^1, f_{j_2}^2, \dots, f_{j_{N_f}}^{N_f}$ with the constraint that the leaks should
301 appear at sequential time instants $k_1 < k_2 < \dots < k_{N_f}$ that allow the leak detection, estimation
302 and localization method to sequentially update the hydraulic simulator with extra demands at leak
303 localization nodes $\hat{j}_1, \hat{j}_2, \dots, \hat{j}_{N_f}$. Once a leak has been fixed, the simulator is updated, eliminating
304 the extra demand related to the fixed leak. This process can be done manually or automatically
305 detecting negative values in the leak estimation, computed by Equation (5).

306 The model-based leak localization method performance depends on the accuracy of the hy-
307 draulic model, sensor noises and the availability of reliable users demand information (Blesa and
308 Pérez 2018). This third factor is potentially the most critical one because it is not easy to estimate
309 user demands with high accuracy if there are no AMRs installed in some network nodes.

310 **CASE STUDY AND RESULTS**

311 The presented methodologies are tested by means of their application to achieve a solution for
312 the BattLeDIM 2020 challenge. This competition aims at objectively comparing the performance
313 of leak detection and isolation approaches, and hence a common benchmark was prepared for all

314 the competitors: L-Town. It is a small and hypothetical town whose WDN is composed by 782
315 inner nodes, 2 reservoirs, 1 tank and 909 pipes. It is represented in Fig. 1. The inner nodes are
316 coloured to indicate their elevation (which is displayed by means of a colour bar). Pressure sensors,
317 reservoirs and tanks are indicated with special markers. The leaks at the provided datasets are
318 highlighted at the corresponding pipes.

319 The network is composed of three different zones:

- 320 • *Area A*: it is the larger area, composed of the nodes with an elevation (see Fig. 1) between 16
321 and 48 m (655 in total). There is a high density of pressure sensors (29 in total) distributed
322 through the area. Besides, the two water inlets to the network are located within this zone,
323 as well as a tank that is filled from this area to provide water to *Area C*.
- 324 • *Area B*: it comprises the nodes that are elevated less than 16 m (31 in total). This zone is
325 connected to *Area A* by a Pressure Reduction Valve or PRV (there are also installed PRVs
326 downstream of the reservoirs of *Area A*) and there is only one pressure meter in the area.
- 327 • *Area C*: it is composed of the nodes elevated over 48 m (92 in total), and it is supplied with
328 water by the previously mentioned tank. There are 3 pressure sensors installed through the
329 zone, as well as 82 Automated Metered Readings (AMRs) that provide demand information.

330 The competition consists of detecting and localizing the maximum number of leaks from the
331 ones occurred during 2019, as well as the ones remaining from 2018 (represented in Fig. 1), by
332 means of readings of pressure sensors, AMRs of *Area C* and flow meters at the output of the
333 reservoirs and tanks. Moreover, a nominal hydraulic model of the WDN is given, although it is
334 affected by inaccuracies at demand patterns, pipe characteristics, valves status, etc.

335 The different features of the three areas composing the network motivated the usage of the
336 proposed methodologies due to their complementary nature. In this case, the model-based approach
337 is utilized over *Area B* and *Area C*: the former has an unique pressure sensor, so the data-driven
338 cannot perform adequately; while the latter has AMRs installed at a high percentage of the nodes,
339 and hence there is accurate demand information that allows to obtain precise results with the model-

340 based technique. Meanwhile, the data-driven method is exploited in *Area A*, due to the high density
341 of installed pressure sensors, as well as the lack of accurate demand information.

342 A schematic flow diagram of the application of each methodology to the benchmark is showed
343 in Fig. 2.

344 To remark that, for their usage in a real-world application, the leak detection and localization
345 stages would be separated, i.e., first the leak would be detected, and then the corresponding
346 localization algorithm would be applied, finding the leak and fixing it. The presented diagrams
347 reflect the necessity of dealing with multiple and simultaneous leaks at the BattLeDIM2020 case.

348 Besides, the localization methods yield a node/group of nodes as the leak candidate/s, so the best
349 candidates, regarding the corresponding criteria for each method, were used to compute the leaky
350 pipes, considering that the leak must be located at the pipe that connects the best two candidates.

351 **Leak detection and estimation**

352 In order to apply the proposed leak detection and estimated method, the network was split into
353 two distinct areas where the demand forecast could be adapted: *Area A* and *Area B*, containing flow
354 meters at the outlet of reservoirs and tanks; and *Area C*, containing AMR devices.

355 *Area A and Area B*

356 Before applying the leak localization methodology, certain data analyses have to be performed to
357 create leak-free historical data using the information from the flow meters during 2018. Measured
358 data are available every 10 minutes, but it was filtered every hour in order to obtain an hourly
359 demand prediction ($T_s = 1$ hour). In addition, daily periodicity has been considered, i.e., $T = 24$.
360 A polynomial has calibrated every hourly flow prediction (feature) variation throughout the year
361 to obtain the demand forecast with a higher accuracy, considering the variation in the behaviour
362 during the year, which could reach more than 5 l/s on the time of the day analysing:

$$363 \hat{y}_{h,q_{day}} = \sum_{i=0}^{n_p} c_i^h q_{day}^i \quad (16)$$

364 where $\hat{y}_{h,q_{day}}$ is the demand estimation at hour $h = 1, \dots, 24$ (first 1 am and the last 00 am) of day

365 q_{day} ; c_i^h are the coefficients of polynomial at hour h and n_p is the order of the polynomial.

366 The leak detection was carried out according to the proposed leak detection and estimation
367 method, considering $N_W = 24$ (i.e. one day), using the history of free leak-data to calculate the
368 maximum possible error and creating a threshold $\lambda = 2.5$ l/s.

369 As the inlets of Area A and B are the same, when a leak is detected considering demand
370 estimation (16), if a drop of pressure is observed in the inner sensor of Area B, it is assumed that
371 the leak is in Area B and otherwise it is assumed that the leak is in Area A.

372 The result obtained in 2018 with all reported leaks that were detected is shown in Fig. 3(a).

373 *Area C*

374 For this area, due to the AMRs measurement corresponding to 89% of the residences, a more
375 accurate forecast demand was calibrated. Using the first week of the inlet flow and the measurements
376 of the AMRs, a constant K was computed, being the percentage value between both flows. The
377 following equation shows the demand forecast for this area:

$$378 \hat{y}(k) = K \sum_{i=1}^{n_m} \mathbf{q}_{AMR_i}(k) \quad (17)$$

379 where $\mathbf{q}_{AMR_i}(k)$ $n = 1, \dots, n_m$ are the flow measurements at instant k of the n_m AMRs installed in
380 Area C. The leak detection method was applied considering $T_s = 1$ h, $T = 24$ and $N_W = 24$, as in
381 Areas A and B, and a estimated threshold $\lambda = 0.3$ l/s, that is much lower than the leaks present in
382 Area C. Using the proposed method furnished the results shown in Fig. 3(b): a leak of 2 l/s at the
383 beginning of the year, which was not reported by the water utility; and another leak at the beginning
384 of July with a magnitude of 4.5 l/s and fixed in August.

385 **Application of the data-driven methodology**

386 The data-driven approach is divided into its two composing stages, in order to properly expose
387 both their separated and complementary operation in the BattLeDIM case study. To remark that, to
388 face the challenge, the method is applied individually to each time-instant, obtaining a localization
389 result for each one of them. The leak detection method provides the leaks starting time, so that the

390 localization can be applied at the moment of their appearance.

391 *Graph-based state interpolation*

392 The graph-based state interpolation is applied first to recover the complete state of the WDN
393 for both nominal and leaky situations. The supplied EPANET (Rossman 2000) model of the WDN
394 is used to extract the network topology, as well as the node elevations to compute hydraulic heads
395 from the pressure measurements. To perform all the necessary EPANET operations, the EPANET-
396 Matlab-Toolkit (Eliades et al. 2016) was employed. Interpolation results are compared in Fig. 4
397 for the case of the leak at pipe $p158$, considering three possible scenarios regarding the availability
398 and nature of the represented data: (a) simulated states from the available hydraulic model for a
399 nominal situation; (b) interpolated states from the available measurements for a nominal case; (c)
400 interpolated states from the available measurements during a leak event. The hydraulic head is
401 represented by the node colour: the darkest the node, the lowest the value (all the interpolations are
402 achieved settling the α parameter of Equation (14) to 1000).

403 The great similarity between the results of a leak-free EPANET simulation (Fig. 4a) and the
404 estimation of the complete state of the network using nominal data (Fig. 4b) confirms the suitability
405 of the interpolation method and the introduced relaxation of the relation among neighbouring
406 nodes. Besides, an important reduction in the state value can be found if comparing the area that
407 is circled in red at the figure associated to the leak event (Fig. 4c) with the same area at the other
408 two subfigures. This indicates the presence of the leak at the highlighted area.

409 *Leak candidate selection method*

410 Once the complete network state has been computed for both the nominal (\mathbf{x}_{nom}) and leak
411 (\mathbf{x}_{leak}) scenarios, the leak candidate selection method can be used. These vectors are arranged to
412 generate a set of n two-dimensional points, represented as showed in Fig. 5a-b. In this case, Fig. 5a
413 is produced using two (different) nominal vectors, in order to compare two healthy-situations, while
414 Fig. 5b considers \mathbf{x}_{leak} to contain the interpolated state in Fig. 4c, which is caused by a leak.

415 On the one hand, Fig. 5a confirms that the cloud of points for two healthy events is mostly
416 aligned, as the slight differences between expected (by the best-fitting line) and actual coordinates

417 of some points are due to the boundary conditions. On the other hand, Fig. 5b shows the leak
418 effect: the cloud of points is no longer aligned, and a set of them is substantially further to the
419 line in comparison with the rest. The points below the line are the most interesting ones to be
420 considered as candidates by the method. Considering the standard deviation of the distance vector
421 to be a threshold for the selection of the candidates, the result depicted at Fig. 5c can be obtained.
422 It includes a representation of the set of candidates in comparison with the complete network (left),
423 and a zoomed subplot highlighting them in colours depending on the probability of being the leak
424 origin (right). The localization is successful due to the reduced pipe distance between the real leak
425 and the most probable pipe, as highlighted at zoomed subplot of the figure through the colour map.

426 *Complete data-driven strategy*

427 Both the graph-based state interpolation technique and the leak candidate selection method are
428 encapsulated into a complete methodology that receives hydraulic measurements from a single
429 time instant and uses them, together with the network topology, to yield a set of candidates to be
430 the leak location. As the datasets of measurements for the complete years are available, the method
431 can be sequentially applied to every time instant (or a subset of them), so that the evolution of the
432 set of candidates can be analysed to assess the existence of leaks.

433 Besides, as there is not a provided nominal dataset, two options are available to compute this
434 information: to use the hydraulic model and to use data from the provided dataset. The former
435 is discarded due to the data-driven nature of the methodology. Therefore, a certain time window
436 must be selected as the nominal reference for the leak localization process. This characteristic can
437 be exploited to obtain more precise localizations in the presence of old leaks, as their effect will be
438 present at both the nominal reference and the leaky data.

439 **Application of the model-based methodology**

440 The model-based method was applied in two different parts of the network comparing the hourly
441 average value of inner pressure sensor measurements with pressure estimations computed every
442 hour by the hydraulic model. Firstly, it was used in *Area C*, using the measurements of the AMRs
443 in the respective nodes on the demands, and an average was implemented in the nodes without

444 AMRs. Secondly, the method was used in *Area B*, using the first-week node demands collected
445 from the provided EPANET model.

446 *Area C*

447 The proposed leak detection and estimation technique is first applied in order to detect leaks and
448 their magnitudes. As explained in the methodology, individual analysis of each leak is needed to
449 discover the probability of the location in the network, always starting with the first leak detected.

450 Fig. 6a-b shows the result of the leakage location in 2018, with the red line indicating the pipe
451 containing the leak. The leaks were analyzed separately, starting with the first leak, and later it was
452 added to the system to study the second leak.

453 Following the analysis of 2019, the leak of Fig. 6a must be added, as it was not fixed in the
454 previous year. Fig. 6c-d shows the result for the localization leaks, remark that the area in yellow
455 is the nodes with more leak potential, all close to the faulty pipe.

456 *Area B*

457 Additionally, the leak detection and estimation method has used the information of the inner
458 pressure sensor installed in this area with the hydraulic model to determine what is the most
459 probable leak node location (see Fig. 3a).

460 A leak was detected each year so the model-based method was applied. Fig. 6e-f show the leak
461 localization results for years 2018 and 2019 respectively.

462 **BATTLEDIM2020 RESULTS**

463 The application of the detection and localization methods to the BattLeDIM 2019 dataset
464 yielded the results at Table 1. They translated into a third place in the competition, with a true
465 positive rate of 43.47% and only 1 false positive, producing savings of 210,772 €.

466 Several facts about these results must be highlighted to explain the methods performance:

- 467 • The model-based approach is applied to *Area B* and *Area C*, that is, to leaks *p280*, *p277*
468 and *p680*. The localization is accurate, with only 1-2 pipes between the real leak and the

469 detected one (the leak from 2018 affecting *Area C*, i.e., *p257*, was detected and localized,
470 but this was not submitted as a result due to a misunderstanding).

- 471 • Regarding *Area A*, the data-driven method localized the detected leaks with a significant
472 accuracy too, except in the case of the leak at *p142*, due to the reduced density of sensors in
473 the area of this localization. Moreover, the application of the data-driven methodology to
474 single time instants at a time allows to obtain incredibly fast localization results, with only
475 a few minutes or hours of difference in some cases, e.g., *p523*, *p827*, *p331* and *p142*.
- 476 • Consulting Table 1, there were several leaks that were not detected, and hence, they were
477 not submitted. However, the localization methodologies were indeed able to locate them
478 correctly, as it can be appreciated at Fig. 7 (to remark that the figures for leaks at *p426*, *p455*
479 and *p879* were obtained by means of custom experiments, that is, using as nominal data a
480 week that allow to get rid of remanent leaks). The detection of these leaks was hindered
481 because they were partially masked by other leaks during the detection phase, and their
482 influence in the localization stage was not considered important enough to be submitted,
483 due to the hard penalty on false positives in the BattLeDIM2020 challenge.

484 CONCLUSIONS

485 This article presents the application of two complementary strategies to solve the leak detection
486 and localization problem for the case study proposed by the BattLeDIM challenge. On the one hand,
487 a data-driven technique is developed to locate leaks by means of the available pressure information.
488 This method is applied to areas with a high density of pressure meters, in order to attain the best
489 possible results. On the other hand, a model-based approach is applied at areas where the existence
490 of demand meters allows to use the available measured demand information to highly increase the
491 localization precision, as well as to zones where the pressure meters density is too low to properly
492 utilize the data-driven approach. Therefore, a maximum reduction of the dependency of the model
493 is achieved while yielding sufficiently accurate localization results.

494 The results of the application of the proposed methodologies are presented for the dataset of
495 2019, used to evaluate the localization approaches in the competition. The performance of the

496 methods was evaluated with the third place at the challenge, demonstrating the suitability of the
497 usage of the proposed approaches, due to their natural synergy when exploited in a complemen-
498 tary fashion: adapting their utilization to the studied network characteristics and the information
499 availability (hydraulic model, sensors...), the advantages of each technique are boosted and their
500 drawbacks are diminished, yielding a flexible and powerful leak localization solution.

501 Finally, it is important to highlight the fact that there are differences between the application
502 of the presented methods to the BattLeDIM2020 case and most real-world cases, i.e., the datasets
503 of the former case are provided to be analysed in an offline manner, while the latter case would
504 be based on an online detection/localization scheme. This fact implies some extra differences that
505 should be pinpointed:

- 506 • The multi-leak problem would exist in both cases, although the importance of a solution
507 would be radically higher in the offline case, due to the accumulation of leaks that may not
508 be repaired. In a real exploitation, the localized leaks are repaired to avoid the water loss.
- 509 • The influence of uncertainties, noise and sensor precision would be higher in the real-
510 world case. This could imply the necessity of performing averaging and noise/deletion
511 pre-processes, as well as increasing the sampling rate and/or number of gathered samples.

512 Considering these facts, several future worklines remain open: the application of the method-
513 ologies to real cases, in order to face real-world conditions; a deep analysis of the previously
514 commented differences, focusing on aspects like the effect of noise, sampling rates, sensorization
515 properties (amount, placement and precision); as well as the improvement of the techniques to
516 solve known and new problems regarding the leak detection and localization field.

517 **DATA AVAILABILITY STATEMENT**

- 518 • Some or all data, models, or code generated or used during the study are available in a
519 repository online in accordance with funder data retention policies: hydraulic data and
520 model provided by [Vrachimis et al. \(2020\)](#) (they can be downloaded from <https://>

521 battledim.ucy.ac.cy/?page_id=33).

- 522 • Some or all data, models, or code generated or used during the study is proprietary or
523 confidential in nature and may only be provided with restrictions (e.g. anonymized data):
524 leak detection and localization algorithms codes.

525 **ACKNOWLEDGEMENTS**

526 The authors want to thank the Spanish national project DEOCS (DPI2016-76493-C3-3-R) and
527 L-BEST (Ref. PID2020-115905RB-C21), as well as the Spanish State Research Agency through
528 the María de Maeztu Seal of Excellence to IRI (MDM-2016-0656).

529 **REFERENCES**

- 530 Adedeji, K. B., Hamam, Y., Abe, B. T., and Abu-Mahfouz, A. M. (2017). “Leakage detection and
531 estimation algorithm for loss reduction in water piping networks.” *Water*, 9(10), 773.
- 532 Arifin, B. M. S., Li, Z., Shaha, S. L., Meyer, G. A., and Colin, A. (2018). “A novel data-driven leak
533 detection and localization algorithm using the Kantorovich distance.” *Comput. Chem. Eng.*, 108,
534 300–313.
- 535 Blesa, J. and Pérez, R. (2018). “Modelling uncertainty for leak localization in water networks.”
536 *IFAC-PapersOnLine*, 51(24), 730–735.
- 537 Chan, T., Chin, C., , and Zhong, X. (2018). “Review of current technologies and proposed intelligent
538 methodologies for water distributed network leakage detection.” *IEEE Access*, 6, 78846–78867.
- 539 Donkor, E. A., Mazzuchi, T. A., Soyer, R., and Alan Roberson, J. (2014). “Urban water demand
540 forecasting: review of methods and models.” *Journal of Water Resources Planning and Man-*
541 *agement*, 140(2), 146–159.
- 542 Duan, H., Lee, P., Ghidaoui, M., and Tung, Y. (2011). “Leak detection in complex series pipelines
543 by using the system frequency response method.” *J. Hydraul. Res.*, 49(2), 213–221.
- 544 Eliades, D. G., Kyriakou, M., Vrachimis, S., and Polycarpou, M. M. (2016). “EPANET-MATLAB
545 toolkit: An open-source software for interfacing EPANET with MATLAB.” *Proceedings of the*
546 *14th International Conference on Computing and Control for the Water Industry, CCWI*.

547 Fanner, P., Davis, S., Hoogerwerf, T., Liemberger, R., Sturm, R., and Thornton, J. (2007). *Leakage*
548 *management technologies*. Water Environ. Research Foundation.

549 Ferrandez-Gamot, L., Busson, P., Blesa, J., Tornil-Sin, S., Puig, V., Duviella, E., and Soldevila, A.
550 (2015). “Leak localization in water distribution networks using pressure residuals and classifiers.”
551 *IFAC-PapersOnLine*, 48(21), 220–225.

552 Goulet, J., Coutu, S., and Smith, I. (2013). “Model falsification diagnosis and sensor placement for
553 leak detection in pressurized pipe networks.” *Adv. Eng. Inf.*, 27(2), 261–269.

554 Han, Q., Eguchi, R., Mehrotra, S., , and Venkatasubramanian, N. (2018). “Enabling state estimation
555 for fault identification in water distribution systems under large disasters.” *2018 IEEE 37th*
556 *Symposium on Reliable Distributed Systems (SRDS)*, IEEE, 161–170.

557 Huang, Y., Zheng, F., Kapelan, Z., Savić, D., Duan, H., and Zhang, Q. (2020). “Efficient leak
558 localization in water distribution systems using multistage optimal valve operations and smart
559 demand metering.” *Water Resour. Res.*, e2020WR028285.

560 Kang, J., Park, Y., Lee, J., Wang, S., and Eom, D. (2017). “Novel leakage detection by ensem-
561 ble CNN-SVM and graph-based localization in water distribution systems.” *IEEE Trans. Ind.*
562 *Electron.*, 65(5), 4279–4289.

563 Kapelan, Z., Savić, D., and Walters, G. (2003). “A hybrid inverse transient model for leakage
564 detection and roughness calibration in pipe networks.” *J. Hydraul. Res.*, 41(5), 481–492.

565 Kim, Y., Lee, S., Park, T., Lee, G., Suh, J., and Lee, J. (2016). “Robust leak detection and its
566 localization using interval estimation for water distribution network.” *Comput. Chem. Eng.*, 92,
567 1–17.

568 LeChevallier, M., Gullick, R., Karim, M., Friedman, M., and Funk, J. (2003). “The potential for
569 health risks from intrusion of contaminants into the distribution system from pressure transients.”
570 *J. Water Health*, 1(1), 3–14.

571 Leflaive, X. (2012). “Water Outlook to 2050: The OECD calls for early and strategic action.”
572 *Global Water Forum*.

573 Li, R., Huang, H., Xin, K., and Tao, T. (2015). “A review of methods for burst/leakage detection

574 and location in water distribution systems.” *Water Sci. Tech.: Water Supply*, 15(3), 429–441.

575 Liemberger, R. and Wyatt, A. (2019). “Quantifying the global non-revenue water problem.” *Water*
576 *Supply*, 19(3), 831–837.

577 Malik, O. (2016). “Probabilistic leak detection and quantification using multi-output Gaussian
578 processes.” Ph.D. thesis, University of Southampton,

579 Menapace, A., Avesani, D., Righetti, M., Bellin, A., and Pisaturo, G. (2018). “Uniformly distributed
580 demand EPANET extension.” *Water Resour. Manag.*, 32(6), 2165–2180.

581 Mohar, B., Alavi, Y., Chartrand, G., and Oellermann, O. (1991). “The Laplacian spectrum of
582 graphs.” *Graph theory, combinatorics, and applications*, 2(12), 871–898.

583 Mutikanga, H., Sharma, S., and Vairavamoorthy, K. (2013). “Methods and tools for managing
584 losses in water distribution systems.” *J. Water Resour. Plan. Manag.*, 139(2), 166–174.

585 Pérez, R., Puig, V., Pascual, J., Quevedo, J., Landeros, E., and Peralta, A. (2011). “Methodology
586 for leakage isolation using pressure sensitivity analysis in water distribution networks.” *Control*
587 *Eng. Pract.*, 19(10), 1157–1167.

588 Pérez, R., Sanz, G., Puig, V., Quevedo, J., Cugueró-Escofet, M., Nejjari, F., Meseguer, J., Cem-
589 brano, G., Mirats-Tur, J., and Sarrate, R. (2014). “Leak localization in water networks: a
590 model-based methodology using pressure sensors applied to a real network in Barcelona.” *IEEE*
591 *Control Syst.*, 34(4), 24–36.

592 Pérez-Pérez, E., López-Estrada, F. R., Valencia-Palomo, G., Torres, L., Puig, V., and Mina-Antonio,
593 J. (2021). “Leak diagnosis in pipelines using a combined artificial neural network approach.”
594 *Control Engineering Practice*, 107, 104677.

595 Rajeswaran, A., Narasimhan, S., and Narasimhan, S. (2018). “A graph partitioning algorithm for
596 leak detection in water distribution networks.” *Comp. & Chem. Eng.*, 108, 11–23.

597 Rossman, L. A. (2000). “EPANET 2: Users Manual.

598 Sanz, G., Pérez, R., Kapelan, Z., and Savić, D. (2016). “Leak detection and localization through
599 demand components calibration.” *J. Water Resour. Plan. Manag.*, 142(2), 04015057.

600 Soldevila, A., Blesa, J., Jensen, T., Tornil-Sin, S., Fernández-Cantí, R., , and Puig, V. (2021).

601 “Leak Localization Method for Water-Distribution Networks Using a Data-Driven Model and
602 Dempster-Shafer Reasoning.” *IEEE Trans. Control Syst. Technol.*, 29(3), 937–948.

603 Soldevila, A., Blesa, J., Tornil-Sin, S., Duviella, E., Fernandez-Canti, R. M., and Puig, V. (2016).
604 “Leak localization in water distribution networks using a mixed model-based/data-driven ap-
605 proach.” *Control Engineering Practice*, 55, 162–173.

606 Sophocleous, S., Savić, D., and Kapelan, Z. (2019). “Leak localization in a real water distribution
607 network based on search-space reduction.” *J. Water Resour. Plan. and Manag.*, 145(7), 04019024.

608 Torres, L., Verde, C., and Molina, L. (2021). “Leak diagnosis for pipelines with multiple branches
609 based on model similarity.” *Journal of Process Control*, 99, 41–53.

610 Vítkovský, J., Simpson, A., and Lambert, M. (2000). “Leak detection and calibration using tran-
611 sients and genetic algorithms.” *J. Water Resour. Plan. Manag.*, 126(4), 262–265.

612 Vrachimis, S. G., Eliades, D. G., Taormina, R., Ostfeld, A., Kapelan, Z., Liu, S., Kyriakou, M.,
613 Pavlou, P., Qiu, M., and Polycarpou, M. M. (2020). “BattLeDIM: Battle of the Leakage Detection
614 and Isolation Methods.

615 Wang, X., Waqar, M., Yan, H., Louati, M., Ghidaoui, M., Lee, P., Meniconi, S., Brunone, B.,
616 and Karney, B. (2020). “Pipeline leak localization using matched-field processing incorporating
617 prior information of modeling error.” *Mech. Syst. Signal Process.*, 143, 106849.

618 Wu, Y. and Liu, S. (2017). “A review of data-driven approaches for burst detection in water
619 distribution systems.” *Urban Water J.*, 14(9), 972–983.

620 Xu, Q., Liu, R., Chen, Q., and Li, R. (2014). “Review on water leakage control in distribution
621 networks and the associated environmental benefits.” *J. Environ. Sci.*, 26(5), 955–961.

622 Zhou, X., Tang, Z., Xu, W., Meng, F., Chu, X., Xin, K., and Fu, G. (2019). “Deep learning identifies
623 accurate burst locations in water distribution networks.” *Water Res.*, 166, 115058.

624 Zhu, X., Ghahramani, Z., and Lafferty, J. (2003). “Semi-supervised learning using gaussian fields
625 and harmonic functions.” *Proceedings of the 20th International conference on Machine learning
626 (ICML-03)*, 912–919.

627

List of Tables

628

1 Results of the BattLeDIM 2020 for the *IRI* team 27

TABLE 1. Results of the BattLeDIM 2020 for the *IRI* team

Real leak	Start time	Detected leak	Detection time	Distance (m)	Distance (n ^o pipes)
p257	2019-01-01 00:05	-	-	-	-
p427	2019-01-01 00:05	-	-	-	-
p810	2019-01-01 00:05	p798	2019-10-25 06:00	237.48	4
p654	2019-01-01 00:05	p662	2019-02-22 07:00	299.33	5
p523	2019-01-15 23:00	p500	2019-01-15 23:10	192.88	3
p827	2019-01-24 18:30	p64*	2019-01-24 18:50	335.93	6
p280	2019-02-10 13:05	p278	2019-02-10 20:00	98.41	1
p653	2019-03-03 13:10	-	-	-	-
p710	2019-03-24 14:15	-	-	-	-
p514	2019-04-02 20:40	p91	2019-04-02 21:00	249.28	5
p331	2019-04-20 10:10	p360	2019-04-20 10:15	278.91	5
p193	2019-05-19 10:40	-	-	-	-
p277	2019-05-30 21:55	p249	2019-06-26 20:00	119.49	1
p142	2019-06-12 19:55	p650*	2019-06-12 19:55	435.27	9
p680	2019-07-10 08:45	p207	2019-07-10 09:00	113.87	2
p586	2019-07-26 14:40	p563	2019-07-31 19:15	224.31	3
p721	2019-08-02 03:00	-	-	-	-
p800	2019-08-16 14:00	p179	2019-08-25 13:25	74.28	1
p123	2019-09-13 20:05	-	-	-	-
p455	2019-10-03 14:00	-	-	-	-
p762	2019-10-09 10:15	-	-	-	-
p426	2019-10-25 13:25	-	-	-	-
p879	2019-11-20 11:55	-	-	-	-

* These leaks are not considered to be detected due to the detection range of 300 m

629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649

List of Figures

- 1 Schematic representation of L-Town. 29
- 2 Schematic flowcharts of the application of the leak detection and localization methodologies: (left) model-based; (right) data-driven. 30
- 3 Graph of leak detection results in 2018 (a) Analysis of the input flow corresponding to Area A and B, with the black markers “o” being the time of detection and the red markers “x” are the leak fix; (b) Analysis of the input flow corresponding to Area C. 31
- 4 Graphical comparison of the interpolated states for the case of a leak at pipe *p158* among the three possible scenarios regarding the availability and nature of the represented data: (a) Nominal EPANET data; (b) Nominal interpolated data; (c) Interpolated data for leak *p158* 32
- 5 Graphical results of the leak candidate selection method: (a-b) Representation of the generated clouds of points for the leak-free (a) and leak (b) scenarios (blue markers) together with the best fitting line (red line); (c) Global localization result showing the complete network and highlighting the leak candidate nodes in yellow (left), and local localization result, illustrated by a colour map with blue representing the least probable candidates, and yellow indicating the most probable ones (right). . . 33
- 6 Graphical representation of the localization result for the following leaks: (a) *p257*; (b) *p31*; (c) *p280*; (d) *p277*; (e) *p673*; (f) *p680*. 34
- 7 Graphical representation of the localization result for the following undetected leaks: (a) *p653*; (b) *p710*; (c) *p193*; (d) *p721*; (e) *p762*; (f) *p426*; (g) *p455* & *p879*. 35

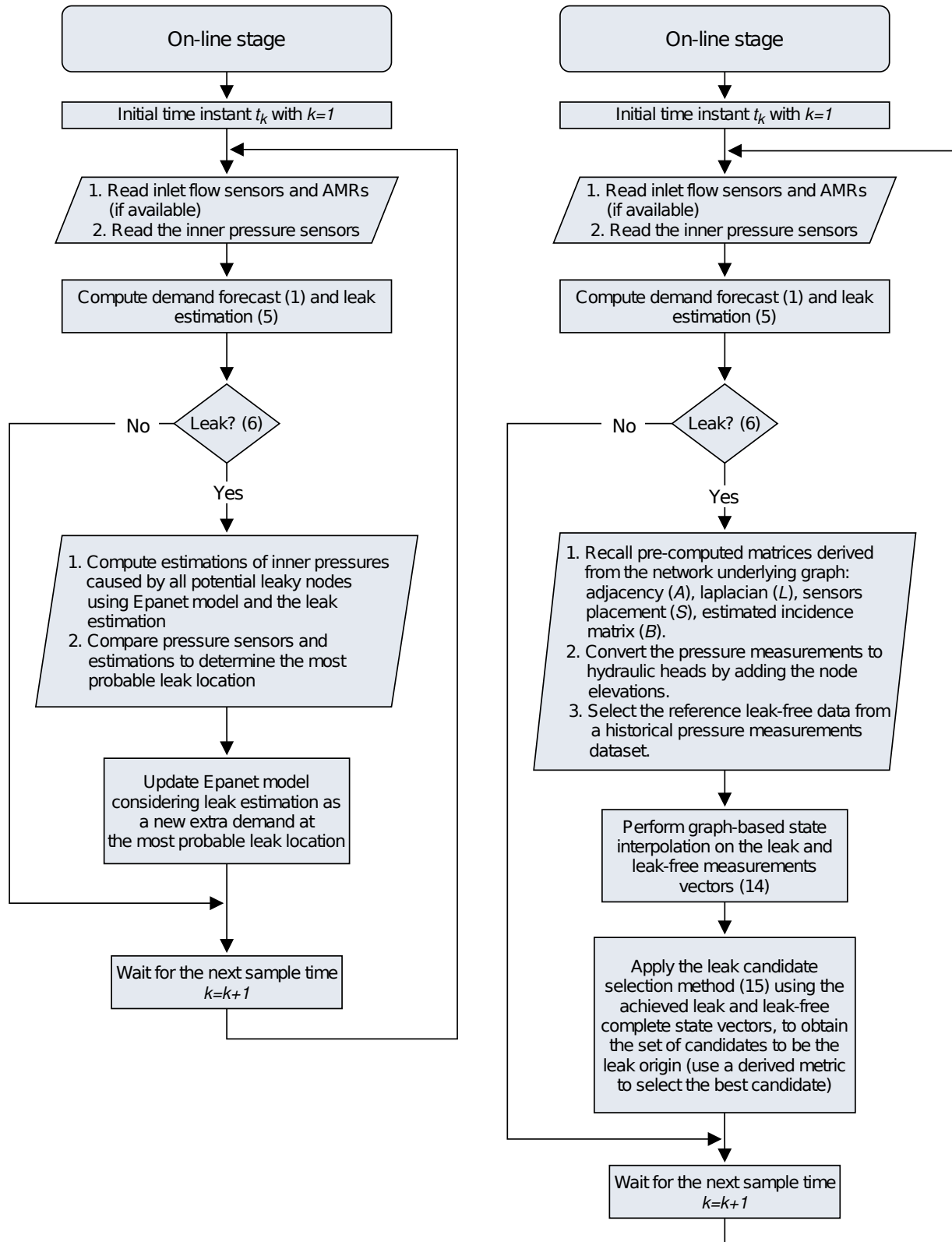
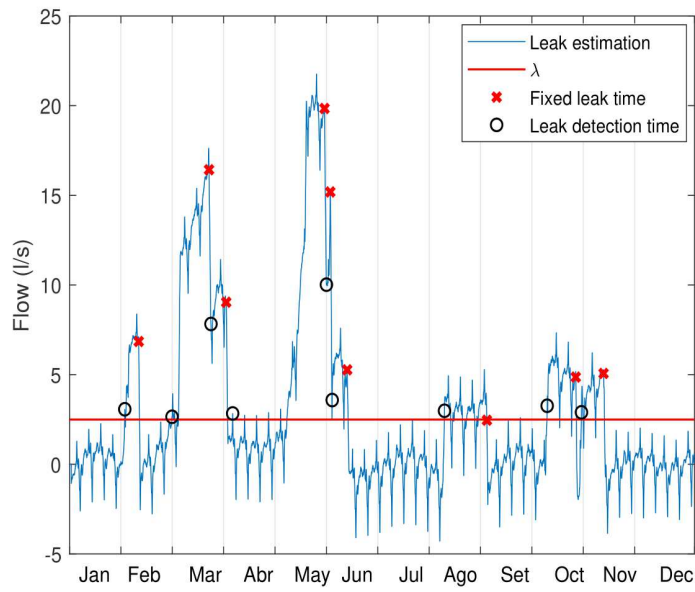
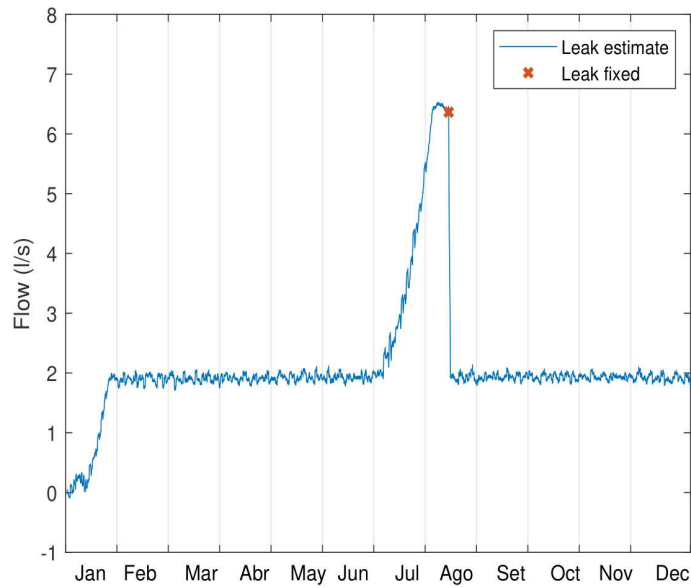


Fig. 2. Schematic flowcharts of the application of the leak detection and localization methodologies: (left) model-based; (right) data-driven.



(a)



(b)

Fig. 3. Graph of leak detection results in 2018 (a) Analysis of the input flow corresponding to Area A and B, with the black markers “o” being the time of detection and the red markers “x” are the leak fix; (b) Analysis of the input flow corresponding to Area C.

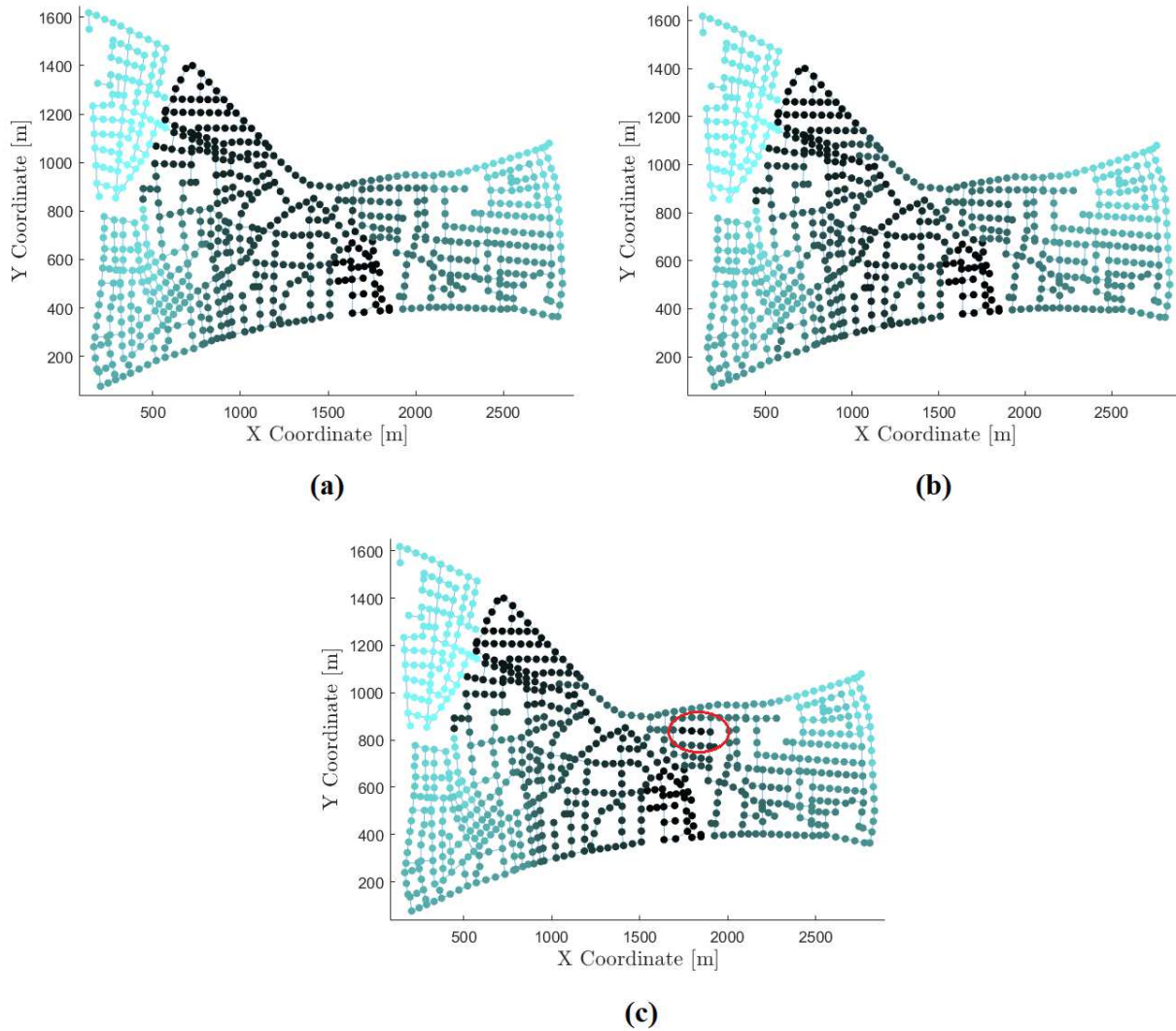


Fig. 4. Graphical comparison of the interpolated states for the case of a leak at pipe $p158$ among the three possible scenarios regarding the availability and nature of the represented data: (a) Nominal EPANET data; (b) Nominal interpolated data; (c) Interpolated data for leak $p158$

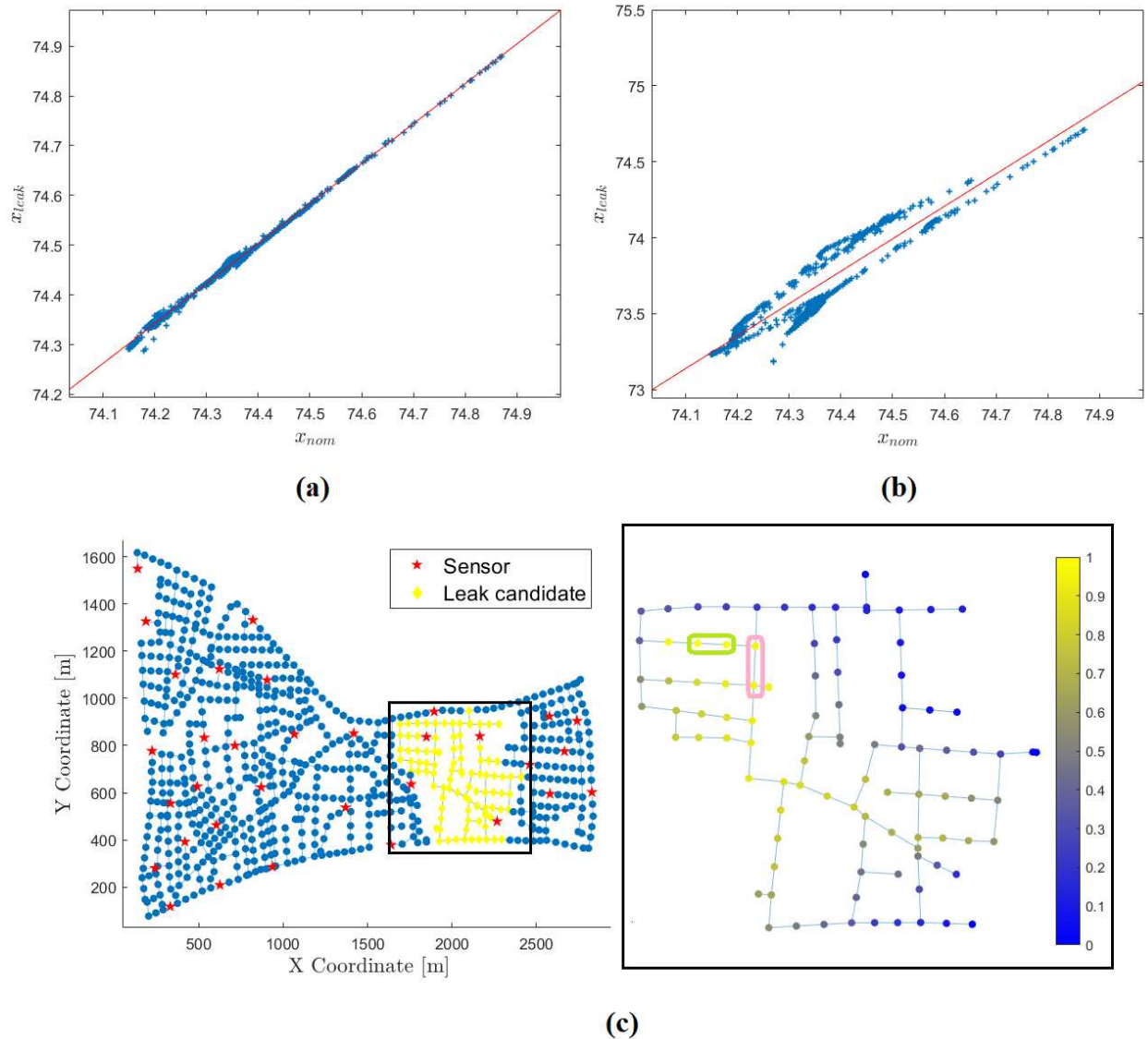


Fig. 5. Graphical results of the leak candidate selection method: (a-b) Representation of the generated clouds of points for the leak-free (a) and leak (b) scenarios (blue markers) together with the best fitting line (red line); (c) Global localization result showing the complete network and highlighting the leak candidate nodes in yellow (left), and local localization result, illustrated by a colour map with blue representing the least probable candidates, and yellow indicating the most probable ones (right).

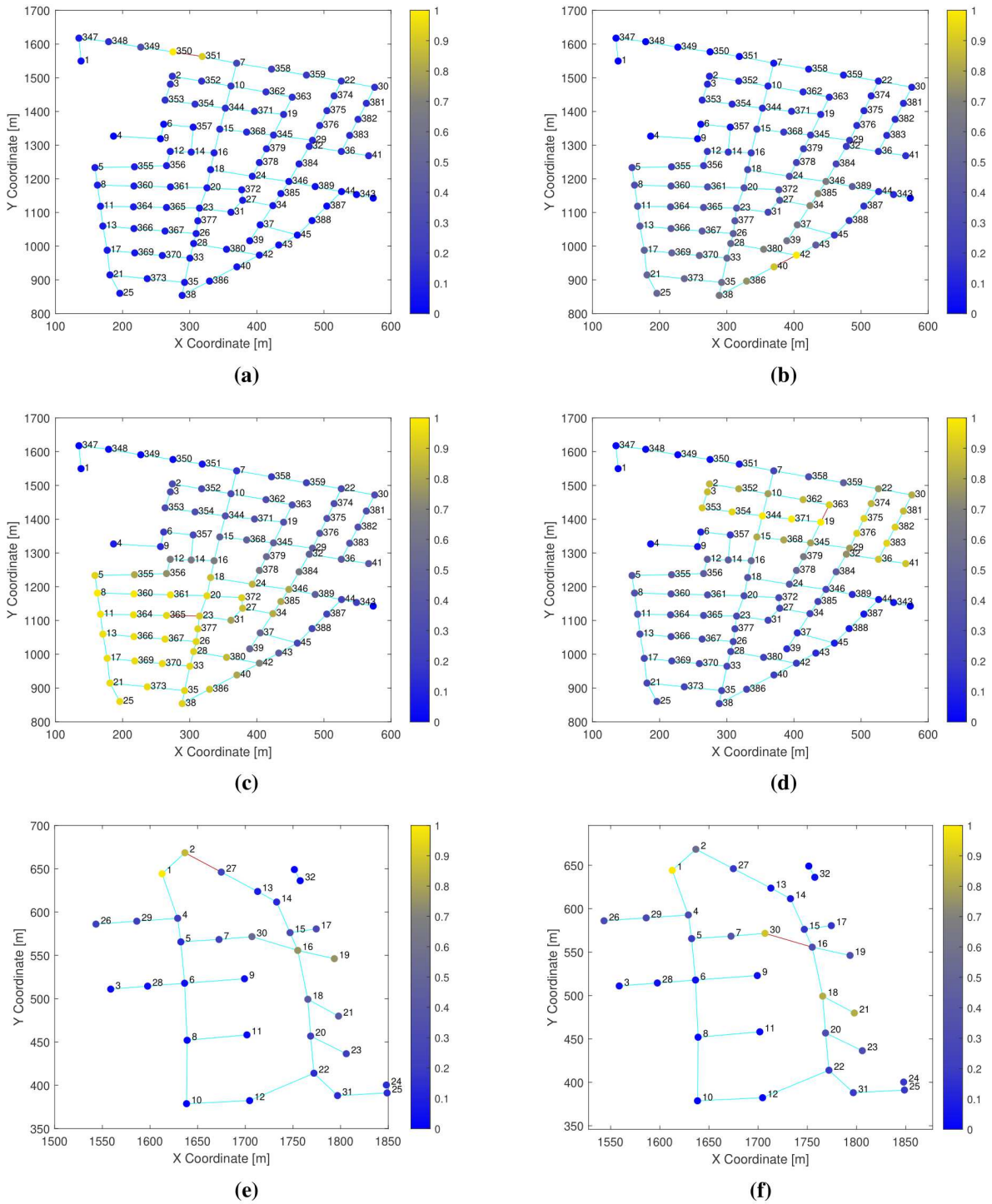


Fig. 6. Graphical representation of the localization result for the following leaks: (a) p_{257} ; (b) p_{31} ; (c) p_{280} ; (d) p_{277} ; (e) p_{673} ; (f) p_{680} .

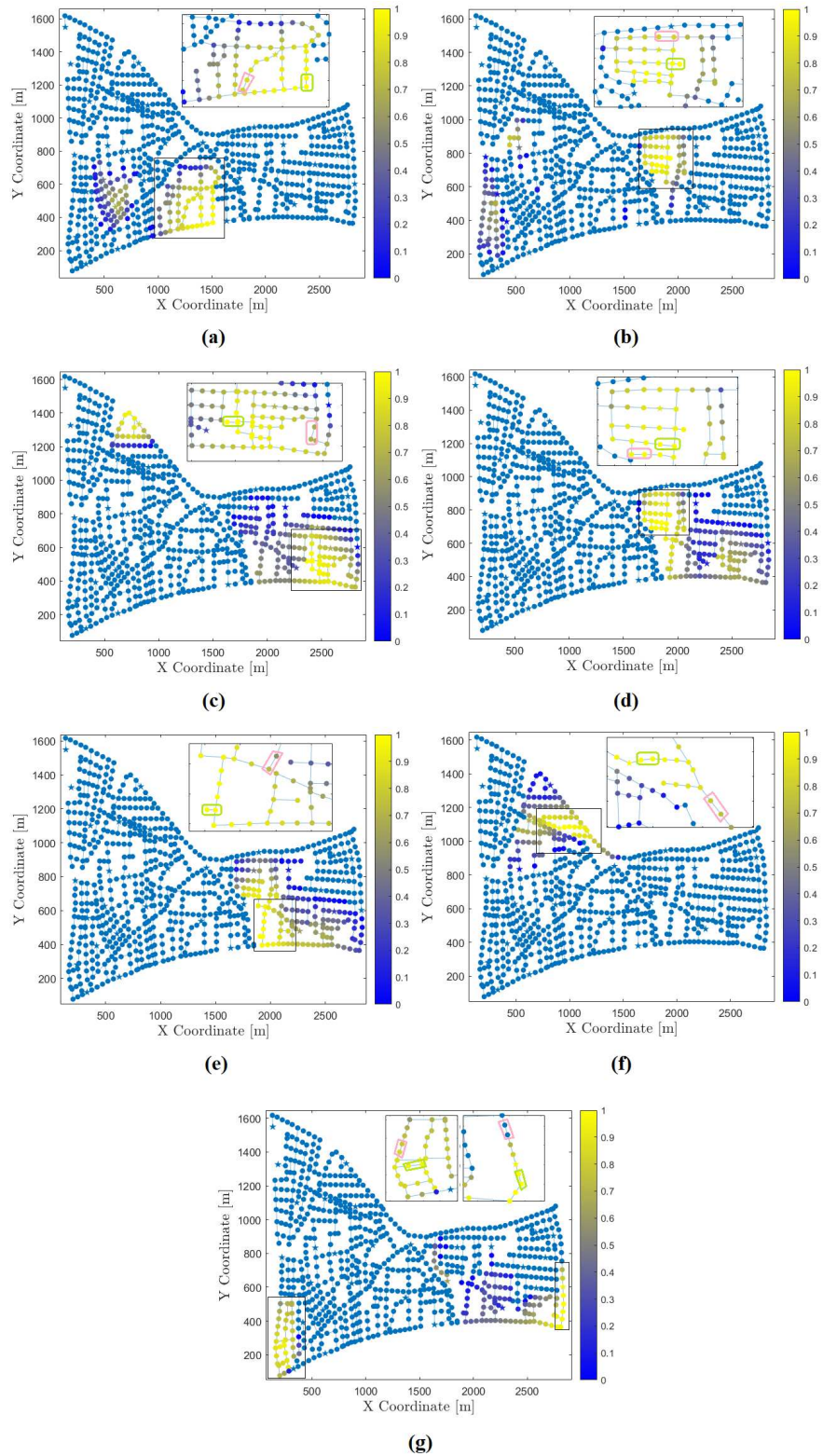


Fig. 7. Graphical representation of the localization result for the following undetected leaks: (a) *p653*; (b) *p710*; (c) *p193*; (d) *p721*; (e) *p762*; (f) *p426*; (g) *p455* & *p879*.