# An effective and efficient method for identification of contamination sources in water distribution systems based on manual grab-sampling

**Yiran Ji[1], Feifei Zheng[2], Jiawen Du[3], Yuan Huang[4], Weiwei Bi[5], Huan-Feng Duan[6], Dragan Savic[7] and Zoran Kapelan[8]**

**[1]Yiran Ji:** PhD student, College of Civil Engineering and Architecture, Zhejiang University, China. yiranji@zju.edu.cn

**[2]Feifei Zheng:** Professor, College of Civil Engineering and Architecture, Zhejiang University, China. feifeizheng@zju.edu.cn. Corresponding author, Tel: +86-571-8820-6757. Postal address: A501, Anzhong Building, Zijingang Campus, Zhejiang University, 866 Yuhangtang Rd, Hangzhou, 310058 China.

**[3]Jiawen Du:** Master student, College of Civil Engineering and Architecture, Zhejiang University, China. dujiawen@zju.edu.cn

**[4]Yuan Huang:** Associate Professor, College of Water Conservancy & Hydropower Engineering, Hohai University, Nanjing, China. huangyuan@hhu.edu.cn

**[5]Weiwei Bi:** Lecturer, College of Civil Engineering, Zhejiang University of Technology, China. weiweibi@zjut.edu.cn.

**[6]Huan-Feng Duan:** Associate Professor, Department of Civil and Environmental Engineering, The Hong Kong Polytechnic University, Hung Hom, Kowloon, 999077, Hong Kong. hf.duan@polyu.edu.hk.

**[7]Dragan Savic:** Chief Executive Officer, KWR Water Research Institute, Dragan.Savic@kwrwater.nl; Professor, Centre for Water Systems, University of Exeter, North Park Road, Exeter, EX4 4QF, United Kingdom; Distinguished Professor, Universiti Kebangsaan Malaysia.

**[8]Zoran Kapelan:** Professor, Delft University of Technology, Faculty of Civil Engineering and Geosciences, Department of Water Management, Stevinweg 1, 2628 CN Delft, Netherlands. z.kapelan@tudelft.nl.

26  **Abstract**: Most of contamination source localization methods for water distribution systems (WDSs)

27  assumes the availability of accurate water quality models and multi-parameter online sensors, which

28  are often out of reach of many water utilities. To address this, a novel manual grab-sampling method

29  (MGSM) is developed to effectively and efficiently locate continuous contamination sources in a

30  WDS using a dynamic and cyclical sampling strategy. The grab samples are collected at a pre-

31  specified number of hydrants by the corresponding teams followed by laboratory tests. The MGSM

32  optimizes the sampling plan at each cycle by making the probability of contamination source(s) in

33  each sub-network as equal as possible, where sub-networks are determined by the selected hydrants

34  and current flow pipe directions. The CS's size is reduced at each cycle by exploting sample testing

35  results obtained in the previous cycle until there are no further hydrants to sample from. Two real-

36  world WDSs are used to demonstrate the effectiveness of the proposed MGSM. The results obtained

37  show that the MGSM can significantly reduce the spatial range of the CS (to about 5% of the entire

38  WDS) for a range of scenarios including multiple contamination sources and pipe flow direction

39  changes. We found that an optimal number of sampling teams exists for a given WDS, representing

40  a balanced trade-off between detection efficiency and sampling/testing budgets. Due to its relative

41  simplicity the proposed MGSM can be used in engineering practice straightaway and it represents a

42  viable alternative to the methods associated with water quality models and sensors.

43  **Keywords**: Water distribution systems, manual grab-sampling method, contamination sources,

44  water quality

45

## 1. Introduction

A water distribution system (WDS) represents a basic lifeline infrastructure that closely relates to the daily life and health safety of its served population (Qi et al., 2018). Typically, a WDS is spatially distributed and thus inherently vulnerable to accidental and/or intentional contamination intrusion (Ostfeld et al., 2014; Yang and Boccelli 2016; Zhang et al., 2020). For instance, over a five-day period in October 2007, a boil-water notice was served on the majority of Oslo, Norway, as a result of a combination of bacteriological, Cryptosporidium oocysts and Giardia cysts found in the samples taken from the WDS (Robertson et al., 2008). More recently, on 26 July 2020, a contamination event was reported in Hangzhou, China, where a sewer pipe was misconnected to a drinking water pipe in a small suburb (ChinaNews, 2020). Unfortunately, these events were not detected by the water quality warning systems of the local water utilities. The events were reported by the residents and/or diagnosed by the hospitals. This implies that monitoring and protecting water quality safety are still nontrivial challenges for many WDSs (Asheri Arnon et al. 2019).

To secure water quality safety in a WDS, extensive studies have been carried out to develop contamination response systems (CRSs) (Giudicianni et al. 2020a). In principle, an effective CRS should at least consist of a contamination warning and source identification (Rodriguez et al. 2021). Regarding the contamination warning, a straightforward manner is to deploy online water quality sensors within the WDS (Hart and Murray 2010). A warning is triggered once the concentration of some particular water quality parameters (e.g., pH, turbidity) is above or below the sensor's safety threshold. Ideally, placing a sensor at each possible location in the WDS can maximise the capability to generate a warning when a contamination intrusion event occurs (Zheng et al. 2018). However, it is difficult, if not impossible, to implement this approach due to the high capital and maintenance costs associated with so many water quality sensors (Winter et al. 2019).

Consequently, many studies have focused on optimally deploying a limited number of water quality sensors to maximize their detection/warning performance (Rathi and Gupta 2014). These studies range from the use of different objective functions to identify appropriate water quality sensor placement strategies (He et al. 2018; Naserizade et al. 2018), to the development of various algorithms to enable effective optimization on this design problem (Hu et al. 2017). More recently, efforts have been increasingly made to identify design solutions that provide a resilient water

75  quality sensor strategy. The approach does not only perform well when all sensors function
76  perfectly, but also can detect contamination events even under possible sensor failures (Ostfeld et
77  al., 2008; Zhang et al., 2020). Typically, the objective functions designed for the water quality
78  sensor placement problems are very complex as different aspects of contamination detection need
79  to be taken into account (e.g., detection likelihood, detection time delay, sensor reliability, different
80  consequences of non-detection, various uncertainties, Khorshidi et al. 2018). Studies have been
81  undertaken to develop various algorithms to effectively identify optimal water quality sensor
82  placement strategies based on these objective functions (Ung et al. 2017). Specifically, those
83  studies focus on developing either sophisticated search algorithms that enhance the design
84  solution's quality (Di Nardo et al. 2018; Hu et al. 2020) or advanced water quality modelling
85  approaches that improve the optimization efficiency (Naserizade et al. 2018; Ohar et al. 2015).

86  In parallel to the research progress on the early warning systems for contamination detection,
87  efforts have also been made to develop various algorithms for sourcing/localizing the
88  contamination injection locations according to the analysis of sensor data (Pries and Ostfeld, 2007).
89  These developments started by using the traditional optimization techniques, such as linear
90  programming (LP) scheme (Pries and Ostfeld, 2006). This was followed by the use of various
91  evolutionary algorithms (EAs) as they possess superior search capabilities compared to the
92  traditional LP and nonlinear programming (NLP) techniques (Pries and Ostfeld, 2008; Hu et al.,
93  2015; Li et al., 2021). While these algorithms have reliable performance in locating contamination
94  sources in hypothetical case studies, their practical application can be highly challenging. This is
95  mainly due to the "equifinality" issue associated with the identification of the source of the incident
96  (Jia et al., 2021a), where many different injection scenarios (contaminant concentration and
97  starting time) indicate a similar contamination impact. To address this issue, the Bayesian based
98  approaches have been proposed to identify contaminant sources, where the location with the
99  highest posterior probability is interpreted as the most plausible (Yang and Boccelli, 2014; Sankary
100 and Ostfeld, 2019; Jerez et al., 2021). More recently, machine learning algorithms have been
101 increasingly employed to facilitate contamination localization, such as the Random Forest
102 algorithm (Grbčić et al., 2020) and Convolutional Neural Network (Sun et al., 2019).

103 Detailed analysis of previous studies in terms of the CRS research shows that the majority of
104 contamination warning and source identification methods rely heavily on an accurate water quality

105 model (Vrachimis et al. 2020). This is one of the main reasons that may hinder their
106 implementation as a well-calibrated water quality model is usually not available for many water
107 utilities (Sankary and Ostfeld 2018). In addition, existing water quality modelling techniques are
108 still incapable of accurately reproducing contaminant reaction dynamics in WDSs, especially for
109 biochemical contaminants (Hart et al. 2019). While online sensors may provide reliable warning
110 information by measuring the contaminant concentration in real-time, they generally can only
111 measure a limited number of water quality parameters such as pH, turbidity, chlorine and
112 conductivity (Sun et al. 2019). Consequently, many other contaminants such as organics and
113 pathogenic microorganisms cannot be detected with certainty using online in-situ sensors. In
114 addition, water quality sensors are often expensive in both the purchase and maintenance,
115 especially for advanced sensors that are used to measure complex substances (He et al., 2018).
116 Therefore, the water quality sensors are often sparsely distributed in many WDSs (Ostfeld et al.,
117 2014).

118 The contamination events within the WDS can be classified into three different types, which are
119 intentional events (Type 1), accidental events (Type 2) and events caused by the WDS itself (Type
120 3). For Type 1, the contamination can be toxic substances that are intentionally injected into the
121 WDS, typically during a short time period. Such events can result in serious consequences and hence
122 need a quick response at all costs (Ostfeld et al., 2014). Type 2 is often represented by the
123 misconnections between water supply pipes and greywater /sewer pipes that have been reported in
124 China (He et al., 2018). Type 3 can be caused by structural damages to pipes (e.g., contamination
125 due to pipe corrosion or leaks, Zhang et al., 2020) or biochemical substances (e.g., microorganisms)
126 activated by the water at a particular level of turbulence (He et al., 2019).

127 Typically, within Types 2 and 3, the contamination exists *continually* in the WDS until the source(s)
128 is localized and eliminated. These contamination substances (e.g., metal, microorganism, organic)
129 often have the following properties: (i) they can be colorless and tasteless, and hence cannot be
130 directly detected by tap-water users; (ii) they do not induce quick, serious public health consequences
131 (i.e., this study focuses on the contamination events with chronic but no acute health effects) and
132 hence their source(s) localization needs to be conducted without interrupting water supply; and (iii)
133 they may not be directly detected by online water quality sensors as the majority sensors typically
134 monitor simple quality parameters such as chlorine, pH, turbidity and conductivity. These properties

5

135    motivate the development of the proposed manual grab-sampling method (MGSM) to efficiently
136    and effectively identify continuous contamination sources of Types 2 and 3 in WDSs.

137    The proposed MGSM is an iterative manual grab-sampling method (MGSM) to enable effective
138    contaminant detection and localization. This is followed by gathering comprehensive water quality
139    parameter information with the aid of laboratory tests. The MGSM is particularly useful for the
140    cases that the online quality sensors are sparsely distributed (or completely unavailable) or sensors
141    cannot measure the contaminants (Wong et al., 2010). The MGSM does not need water quality
142    modelling and can identify the contamination location without encountering the "equifinality"
143    issue. In addition, for the cases that the labour is plentiful with low cost, the MGSM is preferred
144    as it provides the spatial distribution of water quality measurements at a reduced cost when
145    compared to fixed sensors (Mann et al., 2012). Therefore, manual grab-sampling can be an
146    important strategy for water utilities interested in water quality safety in the WDS, which can
147    supplement the information obtained from existing online sensors.

148    Despite the merits and practical significance of the MGSM for the cases with sparsely distributed
149    sensors and relatively low labor costs, relevant research on this topic is surprisingly rare. Amongst
150    few relevant studies, one significant example is from the work of Wong et al. (2010), where a
151    Mixed-Integer Linear Programming formulation is proposed to determine optimal locations for
152    manual grab sampling after a contamination event is detected in a WDS. In their study, the optimal
153    manual grab sample locations are identified by maximizing the total pair-wise distinguishability
154    of candidate contamination events (eliminate unlikely events as much as possible). While Wong
155    et al. (2010) showed that a contamination event can be identified by their proposed method with
156    significantly improved efficiency, its success was conditioned on a few critical assumptions. These
157    assumptions include: (i) each node in the WDS has an equal probability of being the source of
158    contamination intrusion, (ii) only one contamination event can occur in the WDS, and (iii) the pipe
159    flow direction cannot change during the entire sampling process. However, these assumptions can
160    significantly violate the real conditions as the contamination intrusion can occur at any pipe
161    location and a long pipe is typically associated with a higher contamination probability (He et al.,
162    2018). Furthermore, although the probability of simultaneous multiple contamination intrusions is
163    low, their occurrence is still possible in large WDSs (Butera et al., 2021). In addition, flow

164 direction changes are likely to occur in some pipes in a large WDS with multiple supply sources

165 (Qi et al., 2018).

166 The main contribution of this paper is the proposal of an improved water quality MGSM for

167 detecting and localizing continuous contamination sources in WDSs. The newly developed method

168 employs a dynamic and cyclical sampling strategy based on the hydrant locations in a WDS. The

169 novel aspect of the proposed method is the simple and effective way developed to split the network

170 after each round of sampling, thereby significantly enhancing the efficiency of the entire detection

171 process. In addition, the proposed method is novel in that the optimal sampling locations are

172 determined by making the probability of contamination source in each sub-network based on the

173 current flow pipe directions as equal as possible at each cycle. The results of these samples are

174 subsequently analyzed and employed to drive the sampling strategy for the next cycle. It is

175 highlighted that the proposed MGSM is an alternative to these literature methods (sensor-based

176 methods) in the cases where: (a) sensors are sparsely distributed or not available (e.g. lack of

177 existence of suitable sensors), (b) the low-cost labour force is available, and (c) the contamination

178 events have slow or low impacts to the water quality in the WDSs.

179 **2. Methods**

180 The basic premise of the proposed MGSM is: (1) select a given number of sampling points (hydrants

181 of the WDS) in the studied area based on the testing capacity of the laboratory (i.e., the number of

182 samples that can be tested simultaneously) and the number of sampling teams, with all pipes within

183 the candidate area considered as possible contamination sources, (2) narrow down the range of the

184 candidate areas containing contamination source(s) based on sample testing results, and (3) repeat

185 steps (1) and (2) until the range of candidate areas with contamination source(s) cannot be further

186 narrowed down. The key to effectively implementing this new MGSM is how to automatically select

187 the appropriate hydrants in each cycle of the above methodology to reduce the total number of cycles,

188 thereby quickly localizing the pollution source(s) in the WDS. It is noted that every length of pipe

189 between two hydrants within the WDS is considered as the contamination source. Therefore, the

190 proposed MGSM can account for both the scenarios that the contamination sources are in pipes or

191 junctions. While the proposed MGSM is demonstrated using hydrants in this study, any other

192  sampling facilities (e.g., taps) can be easily handled by simply treating them as hydrants within the
193  algorithm implementation.

194  Section 2 presents the details of the proposed MGSM, including the associated theoretical
195  foundations (e.g., the development of the objective function), the MGSM algorithm structure, the
196  illustration of the proposed MGSM and the optimization method to implement the MGSM.

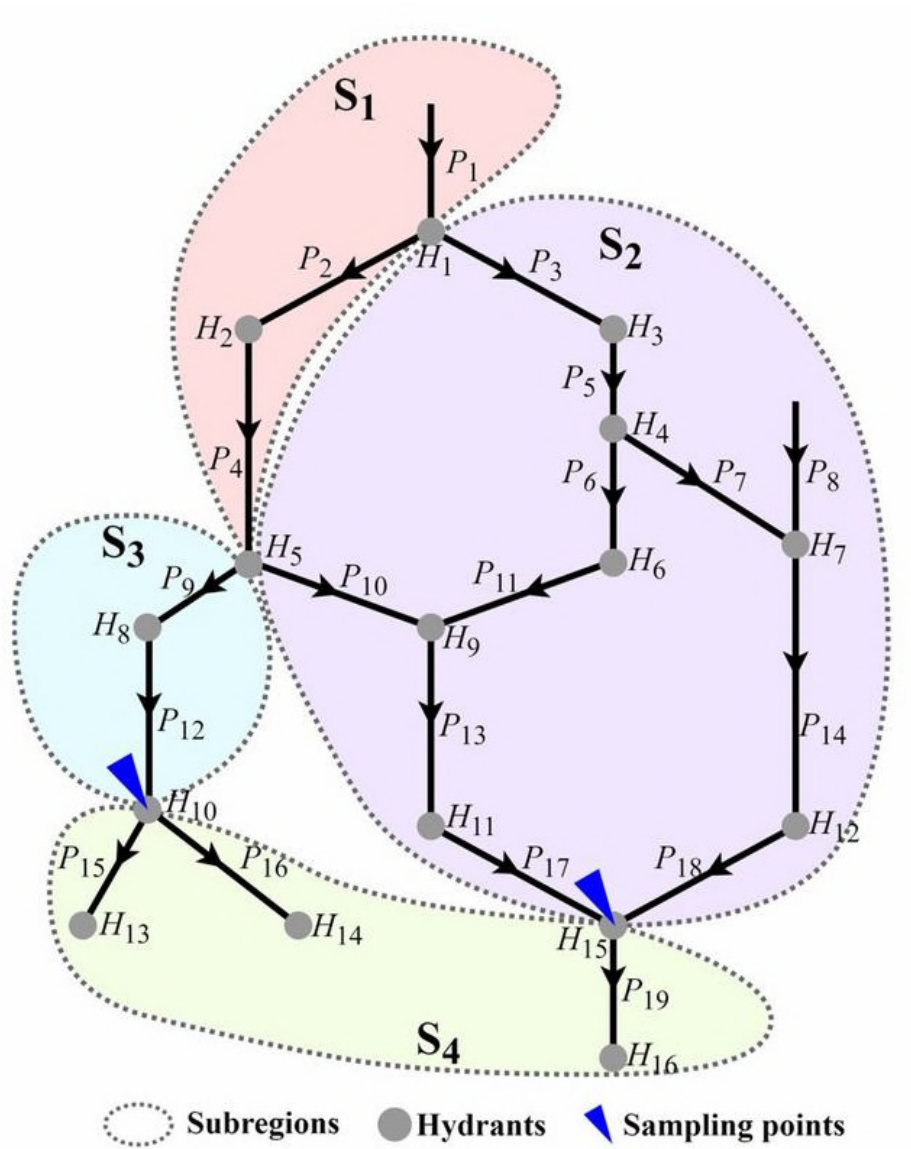**2.1 Theoretical foundations for the proposed MGSM**

197

198  Section 2.1 introduces the theoretical foundations of the proposed MGSM, including the proposal of
199  a method to enable the WDS partitioning and the development of the objective function of the
200  proposed MGSM. The details are given below.

201  **2.1.1 WDS partitioning based on sampling locations and flow directions**

202  As previously stated, the proposed MGSM attempts to identify the optimal sampling locations
203  (hydrants) at each cycle, aimed to minimize the total number of cycles (equivalent to the efficiency
204  and cost of the entire process). Within the MGSM, the entire WDS is partitioned into different sub-
205  networks based on sampling locations and flow directions at a given point in time. Specifically, if a
206  hydrant $H$ in the system is selected as the sampling point, all pipes in the WDS can be divided into
207  two sub-networks: all upstream pipes relative to the selected hydrant $H$, denoted as $\mathbf{U}_H$, and
208  remaining pipes whose flows do not go through $H$, denoted as $\mathbf{N}_H$. If two hydrants ($H_1$ and $H_2$) are
209  selected as the sampling points, four sub-networks can be identified, respectively representing the
210  common group of pipes upstream of both selected hydrants ($\mathbf{U}_1 \cap \mathbf{U}_2$), the unique group upstream of
211  one hydrant only ($\mathbf{U}_1 \cap \mathbf{N}_2$ and $\mathbf{U}_2 \cap \mathbf{N}_1$), and not the upstream of both hydrants ($\mathbf{N}_1 \cap \mathbf{N}_2$). Using this
212  process, for a number of $n$ sampling points in a WDS, e.g., $\{H_1, H_2, \ldots, H_n\}$, a total of $T=2^n$ sub-
213  networks, $\{\mathbf{S}_1, \mathbf{S}_2, \ldots, \mathbf{S}_T\}$, can be obtained theoretically.

214  Figure 1 illustrates how the proposed MGSM identifies the WDS sub-networks based on two
215  sampling locations. A total of 16 hydrants are available that can be considered as the potential
216  sampling points, where the arrows represent pipe flow directions. For illustration, hydrants 10 ($H_{10}$)
217  and 15 ($H_{15}$) are selected as sampling points to enable network partitioning. Four different sub-
218  networks are identified using the proposed MGSM, which are $\mathbf{S}_1=\{P_1, P_2, P_4\}$, $\mathbf{S}_2=\{P_3, P_5, P_6, P_7,$
219  $P_8, P_{10}, P_{11}, P_{13}, P_{14}, P_{17}, P_{18}\}$, $\mathbf{S}_3=\{P_9, P_{12}\}$, $\mathbf{S}_4=\{P_{15}, P_{16}, P_{19}\}$. It can be observed that pipes in $\mathbf{S}_1$

are in the common upstream group for $H_{10}$ and $H_{15}$ and flows for pipes in $S_4$ do not go through any of the two hydrants. Pipes in $S_2$ are those that are upstream of $H_{15}$ but not $H_{10}$, and Pipes in $S_3$ are upstream of $H_{10}$ but not $H_{15}$.



**Figure 1: Illustration of the WDS sub-networks identified by the proposed MGSM based on two sampling locations, with arrows representing pipe flow directions**

For the $n$ sampling points $\mathbf{A}=\{H_1, H_2, \ldots, H_n\}$, the outcome of the test at each sampling point is either that the sample is contaminated or non-contaminated. Therefore, there are $2^n$ possible results for $n$ sampling points, in which each contaminated outcome corresponds to the contamination source

9

229  being located in a certain sub-network or many sub-networks when contaminations are found in
230  many sampling locations. For example, if the contamination is detected at both $H_{10}$ and $H_{15}$, as in
231  Figure 1, it can be derived that the contamination source(s) may be located in the common upstream
232  group of pipes ($\mathbf{S}_1$ in Figure 1). The source can also be in the two sub-networks ($\mathbf{S}_2$ and $\mathbf{S}_3$) upstream
233  of one of the two sampling locations. When only one sampling point indicates contamination, it can
234  be determined that the source is located in the area upstream of the sampling point where
235  contamination is detected, that is, $\mathbf{S}_2$ or $\mathbf{S}_3$. When results show no contamination at both sampling
236  points, then the contamination source(s) is located in an area outside all the upstream parts of the
237  two sampling points, that is, $\mathbf{S}_4$ in Figure 1. This is the basic localization principle used in the
238  proposed MGSM in this study.

239  Once a sub-network or a few sub-networks are selected as potential contamination sources based on
240  the sample testing results, all pipes in this/these sub-network(s) are considered as candidates. This is
241  followed by the further use of the partitioning method to narrow down the spatial range to localize
242  the source. In other words, the network partitioning needs to be carried out at each cycle of the entire
243  sampling process based on the updated candidate pipes with potential contamination sources.

244  **2.1.2 The development of the objective function of the proposed MGSM**

245  Conditioned on the identified $T$ sub-networks, the mathematical expectation ($E(\mathbf{A})$) of a given set of
246  sampling points ($\mathbf{A}$) in localizing the location of the contamination source can be expressed as

$$E(\mathrm{A}) = \sum_{i=1}^{T} p_i \cdot L_i \tag{1}$$

247  where $p_i$ is the probability of the $i^{\text{th}}$ sub-network that have the contamination source, and $L_i$ is the
248  corresponding total pipe length of this sub-network. Since the proposed MGSM mainly aims to
249  detect contamination types 2 and 3 (see section 2 for details), the probability of a contamination
250  source being located on each unit length of pipe can be considered identical. This results in the
251  probability of contamination source being in any sub-network $i$ equal to the ratio of the pipe length
252  of the sub-network $L_i$ to the total pipe length $L_{all}$ in the entire WDS. Mathematically, it gives,

$$E(\mathrm{A}) = \sum_{i=1}^{T} \frac{L_i}{L_{all}} \cdot L_i = \frac{1}{L_{all}} \sum_{i=1}^{T} L_i^2 \tag{2}$$

253    Thus, the objective function for calculating the optimal sampling group can be expressed as follows:

$$\text{Minimize: } F(A) = \frac{E(A)}{L_{all}} = \frac{1}{L_{all}^2} \sum_{i=1}^{T} L_i^2 \tag{3}$$

254    where $F(A)$ is a dimensionless number by dividing $E(\mathbf{A})$ using $L_{all}$, representing the ratio of candidate
255    area with contamination source identified by the sampling group relative to the total pipe length of
256    the entire WDS being considered. $\mathbf{A}$ is the decision variables, representing the hydrant sampling
257    strategy. The minimization of $F(A)$ physically indicates a minimum pipe length of the sub-network
258    with contamination source(s) to be identified by the selected sampling points.

259    Cauchy–Schwarz Inequality (Bhatia and Davis, 1995) can be used to further explain the
260    minimization of Equation (3), which is

$$T \times ( L_1^2 + L_2^2 + \cdots + L_T^2 ) \geq ( L_1 + L_2 + \cdots + L_T )^2 \tag{4}$$

$$\text{Namely } F(A) = \frac{1}{L_{all}^2} \sum_{i=1}^{T} L_i^2 \geq \frac{1}{T} \tag{5}$$

261    For $L_1 = L_2, \ldots, = L_T$, the equation holds. Under this condition, when only one hydrant is selected as the
262    sampling point in each cycle, the optimal hydrant divides the WDS into two sub-networks such the
263    pipe length of its upstream section is half of the total length. When $n$ hydrants are selected as the
264    sampling points in each cycle, theoretically, the optimal hydrant group bisects the WDS to $2^n$ sub-
265    networks with identical pipe lengths across different sub-networks. In other words, the minimization
266    of Equation (3) (i.e., $L_1 = L_2, \ldots, = L_T$) can be interpreted as using a specified number of sampling points
267    to assign the pipes into $T$ sub-networks with the minimum difference in pipe length at each cycle.
268    This is equivalent to the bi-section approach in computer science, and hence it is expected that such
269    a method can achieve a statistically efficient sampling strategy to localize the contamination source.
270    It is noted that the proposed optimization method may not be able to guarantee global optimality, but
271    it can offer a near-optimal solution that can be efficiency found at each cycle.

272    The pipe length is used to split the WDS in this study due to its simplicity and efficiency. However,
273    a more refined method may need to account for water velocities or flow volumes, both of which can
274    be correlated with pipe diameters, as well as can account for the amounts of contaminants moving

11

275 through the pipes. Therefore, partitioning the WDS with the aid of both pipe length and water
276 velocity can be an important future research focus.

## 2.2 The algorithm of the proposed MGSM

278 The implementation of the proposed MSGM can be triggered by (i) the routine water quality
279 checking operation required by the water utilities, (ii) abnormal signals from online water quality
280 sensors (e.g., chlorine sensors) that are often installed at the outlets of the districted metering areas
281 (DMAs), or (iii) positive testing results of samples at the outlets of the DMAs or at the important
282 locations within the WDS area. Figure 2 shows the algorithm details of the proposed MGSM in
283 localizing contamination source(s). As shown in this figure, when the number of sampling locations
284 at each cycle is $n=1$, the sampling hydrant is selected by minimizing Equation (3), where the
285 minimization method is elaborated in Section 2.4. The candidate sub-network (**CS**) that may contain
286 contamination source(s) is updated at each cycle based on the sample testing results (Case A1 and
287 Case A2 in Figure 2). If $n$ is greater than 1, the algorithm of the proposed MGSM becomes more
288 complex, with details given in Figure 2. At the beginning (i.e., flag=0, and the MGSM is triggered),
289 the $n$ optimal sampling locations are identified by minimizing Equation (3) for the entire WDS being
290 considered (i.e., **CS** is the entire WDS). This is followed by the application of selection strategy 1
291 (SA1) to update the **CS** for the next cycle, where three different cases (Case B1, B2 and B3) can be
292 available. For Case B2 (only one sample hydrant has contamination) and B3 (all sample hydrants
293 are contamination free), it is straightforward to select the **CS** for the next cycle as shown in Figure
294 2.

---

**Specify** the number of sampling points $n$
**Set** the cycle $c=1$, flag=0, the candidate sub-network (**CS**) as the entire WDS
**While** True
{
    **If** $n = 1$
    {
        **Select** $n$ sampling hydrant for the **CS** by minimizing Equation (3)
        **Update** the **CS** according to sample testing results
          Case A1: the sample is contaminated
               Select the sub-network (**US**) upstream of the selected hydrant
          Case A2: the sample is contamination free
               Select the sub-network that is not the upstream of the selected hydrant
        $c = c + 1$
    }
    **Else**
    {
        **If** flag=0

```
{
    Select n sampling hydrants for the CS by minimizing Equation (3)
    Update the CS according to sample testing results using Selection strategy 1 (SA1):
        Case B1: more than one hydrant sample are contaminated
                Select the common sub-network (CUS) upstream of the contaminated hydrants
                Set flag=1
        Case B2: only one hydrant sample is contaminated
                Select the unique sub-network upstream of the contaminated hydrant
        Case B3: no hydrant samples are contaminated
                Select the sub-network that is not the upstream of the selected hydrant
    c = c + 1
}
If flag=1
{
    If the CUS exists and its most downstream hydrant is not sampled
            Assign one sample point at the most downstream hydrant of the CUS
            Select n-1 sampling hydrants for the CS by minimizing Equation (3)
            Update the CS(s) according to testing results at the end hydrant
                If the end hydrant is contaminated
                    Select the CS using the SA1 mentioned above
                Else
                    Selection strategy 2 (SA2):
                    Select the CS(s) as the union of USs of the hydrants showing evidence of
                    contamination minus the union of USs of contamination-free hydrants and the CUS
            c = c + 1
    Else If the CUS does not exist
            Select the CS(s) using the SA2 mentioned above
    Set flag=0
    }
}
If no hydrant can be sampled in the selected sub-network
        break
}
```

**Figure 2: The algorithm of the proposed MGSM**

When more than one sample hydrant is contaminated (Case B1), the common upstream sub-network (**CUS**, which is theoretically available) is selected as the **CS** for the next cycle ($c=c+1$). If this **CUS** exists and its most downstream hydrant is not sampled, one sampling location is assigned to this hydrant. The remaining $n$-1 sampling locations are determined by minimizing Equation (3). The **CS**, which is temporally considered as the **CUS**, is now updated using the following method based on test results of the most downstream hydrant. If that hydrant is contaminated, the **SA1** is employed to update the **CS**, otherwise, the **SA2** (see Figure 2) is used to update the **CS**. Specifically, the **SA2** selects the **CS**(s) as the union of all upstream sub-networks (**US**s) of hydrants where contamination

13

304  was detected, minus the union of **US**s of contamination-free hydrants and the **CUS**. Note that if the
305  selected **CUS** does not exist in the WDS, the **SA**2 is used to update the **CS**(s).

306  The proposed MGSM in Figure 2 can handle both the single and multiple contamination sources in
307  a DMA of a WDS. However, each MGSM run identifies only a sub-network that contains a
308  contamination source of the smallest spatial extent. This identified region may need to be blocked
309  for engineering operations (e.g., disconnect the misconnections, repair the leaks, or replace the pipes),
310  to remove the contamination source(s). Sampling tests with a few contaminated hydrants may
311  indicate the presence of multiple contamination sources in different WDS regions. For such cases,
312  once the identified contamination source(s) is fixed, the proposed MGSM can be applied to the
313  potential **CS**s (instead of the entire WDS) derived by the sampling test results combined with
314  knowledge of pipe flow directions. Such a CS selection can be easily performed by engineering
315  experience, but it is difficult to be shown by formal procedures. However, it is also straightforward
316  to simply apply the MGSM to the entire WDS to identify the other contamination source(s), after
317  the already localized source(s) are fixed.

318  The methodology assumes that all hydrants selected in one cycle can be sampled at the same flow
319  direction status. This assumption is practically reasonable as the time required to grab samples is
320  often short and the frequency of flow direction change is typically low (e.g., once a day, Wong et
321  al., 2010). While flow direction changes may exist within the supply boundary of some real large
322  WDSs, its associated region is often rather small. Therefore, the change of the flow directions will
323  not significantly affect the application of the proposed MGSM. If the WDS region with changing
324  flow direction is large and known, it can be easily accounted for by the proposed MGSM based on
325  an important assumption. This assumption is that the time between the start of the flow direction
326  change and the next sampling cycle is significantly greater than the longest travel time from the
327  source to the sample locations. In other words, the contaminant distribution has to be consistent with
328  the current flow regime and can have no residual effects from the previous flow regime. Based on
329  this assumption, the flow direction changes can be considered by the WDS partitioning process as
330  described in Section 3.1.1, which would accordingly affect the formulation of sub-networks and
331  hence the identification of the optimal sampling locations (Equation 3).
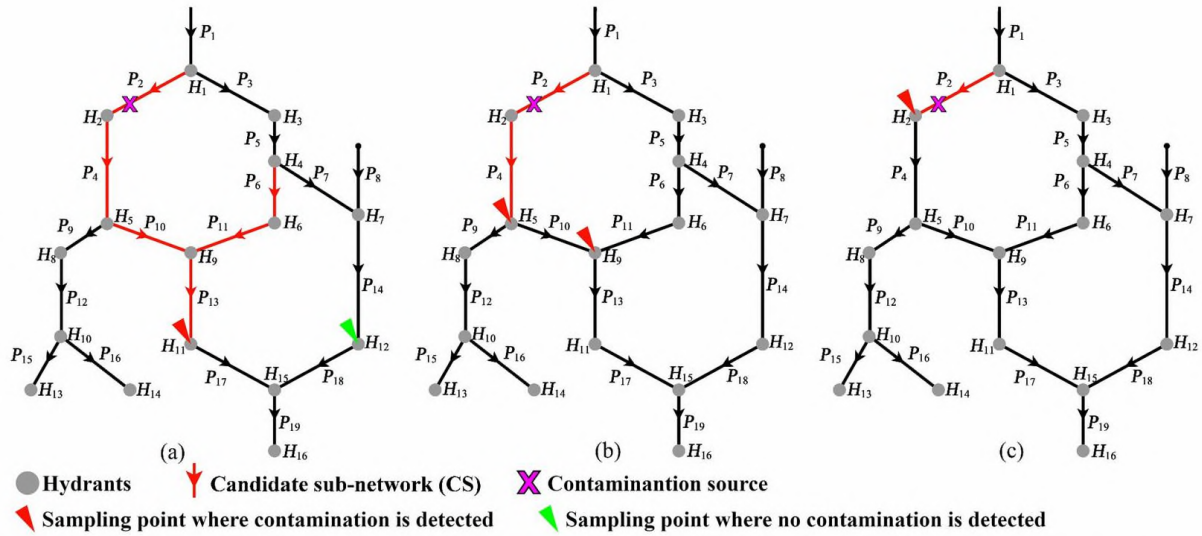
332  **2.3 Illustration of the proposed MGSM**

333  The proposed MGSM is illustrated with two scenarios, including the single contamination source

334  and the two contamination sources simultaneously exist in the WDS, with details given below.

**2.3.1 Single contamination source**

336  We first illustrate the application of the proposed MGSM (Figure 2) using a single contaminating

337  source as shown in Figure 3. The single contamination source is in $P_2$, and two sampling locations

338  ($n$=2) are identified at each cycle. At the first cycle, the entire WDS is set as a candidate sub-network

339  (**CS**), and a total of 120 sampling combinations (two out of 16 total hydrants) are possible. The

340  mathematical expectations (Equation 3) corresponding to these 120 combinations are calculated by

341  enumeration and the combination with the minimum $F(\mathbf{A})$ value is selected. Consequently, two

342  hydrants $\{H_{11}, H_{12}\}$ are identified as the sampling points yielding the lowest objective function value

343  (Equation 3), as shown in Figure 3(a). Based on the assumed location for the contamination source,

344  the sample from hydrant $H_{11}$ is contaminated while the sample from $H_{12}$ is not based on the

345  laboratory tests. Therefore, the **CS** is updated to be a unique sub-network upstream of $H_{11}$ (and not

346  pipes upstream of $H_{12}$) based on Case B2 in Figure 2, that is, the red pipes shown in Figure 3(a).

347  In the second cycle of sampling, the mathematical expectations corresponding to different hydrant

348  groups are calculated according to the updated **CS** determined in the previous cycle. The resultant

349  optimal strategy is the combination of $H_5$ and $H_9$ as it produces the lowest objective function value.

350  Testing results on these two hydrant samples show that both are contaminated, indicating that the

351  contamination source exists in the common upstream sub-network (**CUS**) of $H_5$ and $H_9$. Therefore,

352  the **CS** is updated as the **CUS** based on Case B1 (Figure 2), which is $\{P_2, P_4\}$ as represented by red

353  lines in Figure 3(b). In the third cycle of sampling, there is only one hydrant location, $H_2$, so the

354  contamination source is successfully detected on $P_2$, which is the exact location of the contamination
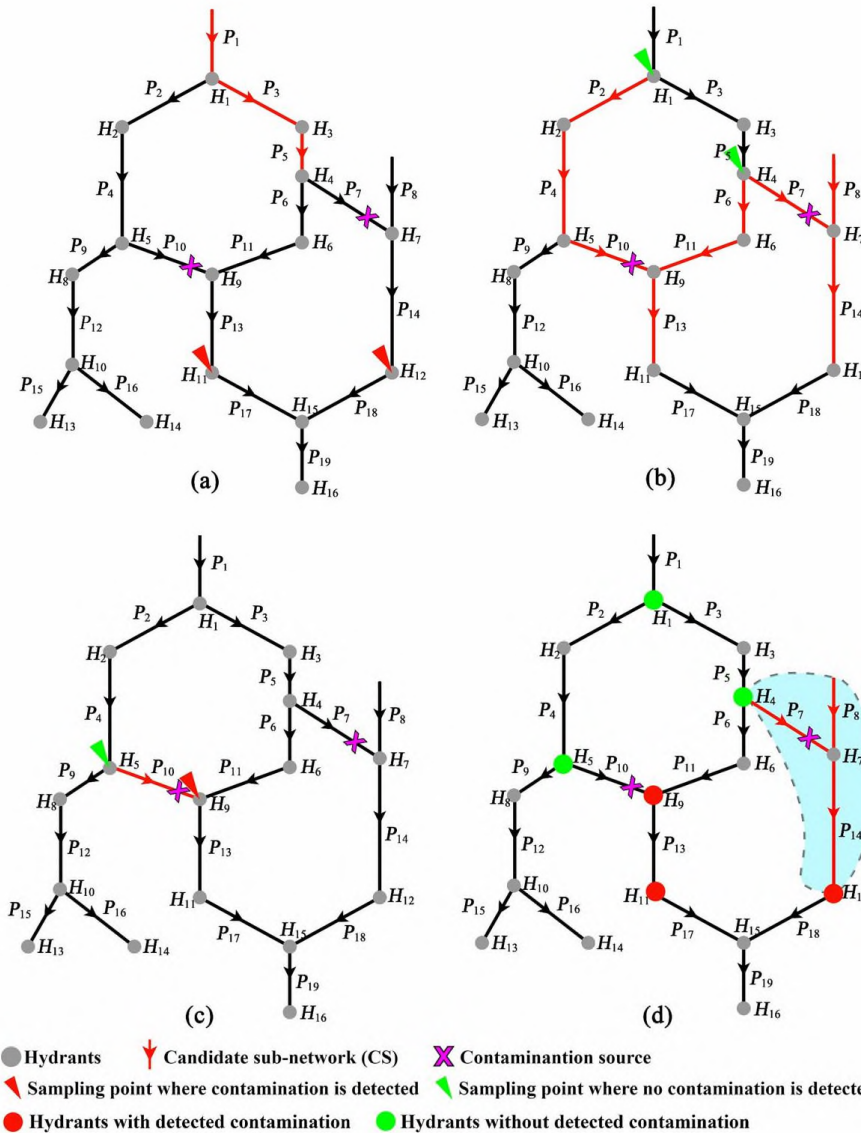
355  source.

**Figure 3: Source localization process for the contamination at $P_2$: (a) the first cycle ($c$=1) of sampling and testing; (b) sampling and testing at $c$=2; (c) sampling and testing at $c$=3**

**2.3.2 Two contamination sources**

Figure 4 illustrate the application of the proposed MGSM (Figure 2) in dealing with two contamination sources. In this figure, the contamination sources are in $P_7$ and $P_{10}$, and two sampling locations ($n$=2) are identified at each cycle. As the same with the single contamination source in Figure 3(a), the hydrants $H_{11}$ and $H_{12}$ are selected as the sampling points at the first cycle by minimizing Equation (3) (the enumeration method is used for this small WDS). The testing results show both hydrants are contaminated, and accordingly, the **CS** is updated to be the common upstream sub-network (**CUS,** red pipes in Figure 4(a)) using Case B1 in Figure 2. Since the **CUS** exists and its most downstream hydrant ($H_4$) is not sampled, $H_4$ is selected as one sampling location and the other location ($H_1$) is identified with the aid of Equation (3) in the second cycle ($c$=2).

Based on the locations of the two contamination sources, the end hydrant $H_4$ should show no contamination in the laboratory test and selection strategy 2 (**SA2**) is used to update the **CS**. More specifically, for such cases, the **CS** can be described as **UA**-**UB**-**CUS** (**CUS** ={$P_1$, $P_3$, $P_5$}), where **UA** is the union of sub-networks (**US**s) upstream of contaminated hydrants (i.e., $H_{11}$ and $H_{12}$ at $c$=1) and **UB** is the union of **US**s sampling hydrants without contaminations (it is null at $c$=1). This is followed by the application of the proposed method at $c$=3, where two hydrants ($H_5$ and $H_9$) are selected as the sampling points. The resultant **CS** is $P_{10}$ using Case B2 in Figure 2 based on test

results ($H_5$ is not contaminated, but $H_9$ is), which is the unique upstream sub-network of $H_9$. Since no hydrants can be sampled in the current **CS** (i.e., $P_{10}$), $P_{10}$ is successfully identified with the contamination source. The run of the proposed MGSM (Figure 2) is finalized.



**Figure 4: Source localization process for two contamination cases at $P_7$ and $P_{10}$: (a) the first cycle ($c$=1) of sampling and testing; (b) sampling and testing at $c$=2; (c) sampling and testing at $c$=3; (d) the CS identified (shaded pipes) for the next MGSM run, where the red and green dots represent test results of the previous MGSM run**

To identify the second contamination source in $P_7$, the s localized source in $P_{10}$ needs to be fixed before the implementation of the next MGSM run. This is because the proposed MGSM identifies

386  only one contamination source for each run. Prior to the application of the next MGSM run, the

387  identified contamination source(s) need to be eliminated. In addition, all the test results of hydrant

388  samples and pipe flow direction information can be jointly used to derive the potential **CS** for the

389  next MGSM run. For the given example, the **CS** can be identified as the red pipes in Figure 4 (d)

390  based on the test results of the previous MGSM run (red and green dots) since (i) the test on $H_4$

391  shows no contamination but $H_{12}$ does, and (ii) the identified source at $P_{10}$ is not upstream of $H_{12}$.

392  This **CS** is only a small proportion of the entire WDS, thereby greatly improving the efficiency of

393  the next MGSM run. However, for cases when the **CS** cannot be determined by the existing

394  information provided by sample test results and pipe flow directions, the entire WDS (after the

395  identified contamination source(s) is eliminated) is considered as the **CS** again to enable the

396  application of the proposed MGSM.

397  In this subsection, one and two contamination sources are used to illustrate the proposed MGSM due

398  to the high likelihood of those events occurring in real WDS. In addition, two sampling locations are

399  used at each cycle for illustration purposes, where the pipe flow directions are not changed. However,

400  the application procedures with details given in Figure 2 are generic, and hence can be applied to

401  other scenarios such as different number of sampling locations, different contamination sources and

402  the WDS with possible pipe flow changes (further explanation of which is given in Section 4)
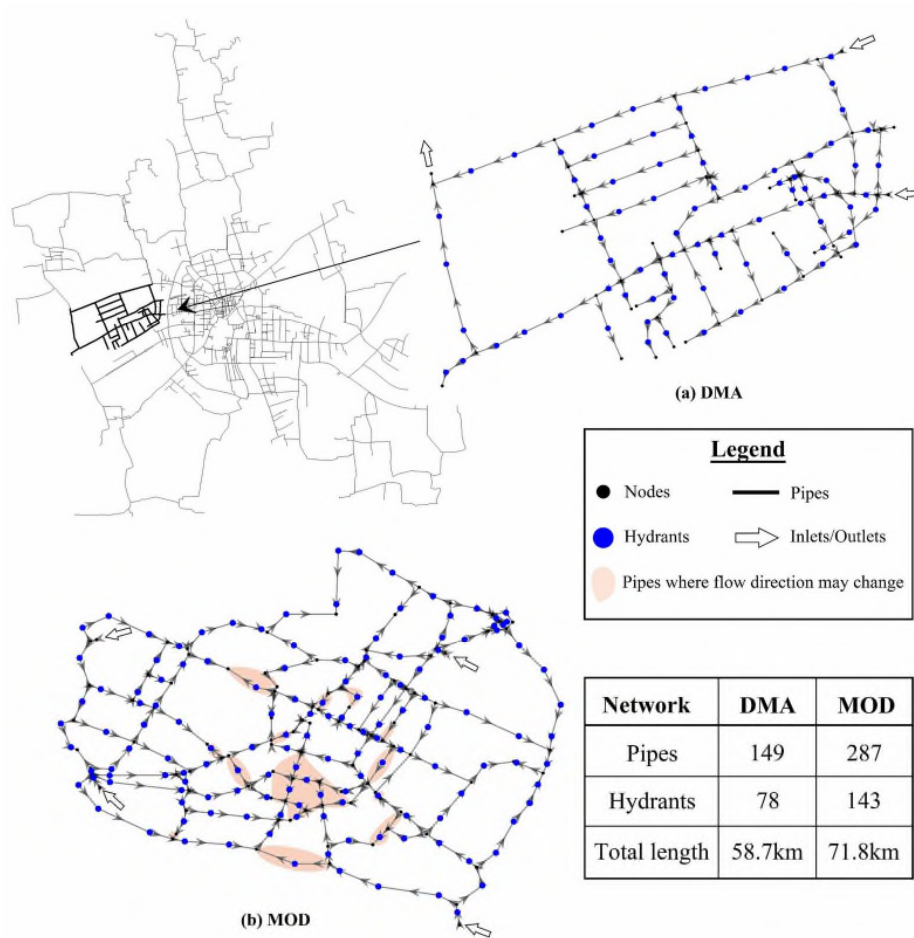
403  **2.4 Optimization method to minimize the objective function**

404  As shown in Figure 2, the proposed MGSM algorithm requires an optimization method to minimize

405  the objective function (Equation 3). While the enumeration method can be effective when dealing

406  with small WDSs and with a low number of sampling locations at each cycle, it is computationally

407  intractable for real and large WDSs. More specifically, for a case with $n$ sampling points applied to

408  a WDS with a total of $N$ hydrants, the number of all possible combinations is $C_N^M$. This value

409  increases exponentially with $n$ and $N$ becoming larger, leading to a rapid increase in computing time

410  and deterioration of detection effectiveness.

411  To solve the computational issue, the Monte Carlo (MC) method is used in this study as an alternative

412  to the enumeration approach in the process of determining the optimal sampling group to improve

413  detection efficiency for large-scale WDSs. The selection of the MC method is mainly due to its

414 simplicity and reasonable performance in offering near-optimal solutions (Maier et al., 2014). This

415 is practically meaningful as in many engineering cases providing near-optimal solutions within a

416 given time framework are more important than identifying global optimums with large

417 computational overheads (Maier et al., 2014). Nevertheless, an advanced optimization algorithm can

418 be developed for the proposed MGSM in future, which is not the focus of the present paper.

419 **3. Case studies**



420

421 **Figure 5: (a) the DMA case study and (b) the MOD case study, where arrows indicate flow**

422 **directions**

423 Two distribution networks (Figure 5) are used to demonstrate the utility of the proposed MGSM.

424 Specifically, the DMA (district meter area) case study is a part of a real-world WDS in China (Figure

425 5a) that consists of 149 pipes (58.7 km in length) and 78 fire hydrants. It has two inlets and one outlet,

19

426    and the flow direction in this network (shown in Figure 5(a)) does not change. The MOD pipe
427    network is a benchmark WDS of the city of Modena in Italy (Bragalli et al., 2012). This network
428    consists of 4 reservoirs (sources), 287 pipes (71.8 km in length) and 143 fire hydrants. Due to the
429    water level changes in the four reservoirs and variations in residential water consumption, the flow
430    directions of some pipes (shaded pipes in Figure 5(b)) in the MOD network change over time.

431    While the demonstration of the proposed MGSM using a very large WDS is academically necessary,
432    but in practice the MGSM is mainly used for a DMA or a region of the entire WDS. This is because
433    (i) many WDSs have been managed by DMAs, which can greatly enhance the operation efficiency,
434    and (ii) for the WDSs with no DMAs, water quality safety checking or contamination sourcing is
435    likely to be conducted region by region. It is highly unlikely to simultaneously consider all the pipes
436    of the entire large network as the contamination sources. Therefore, we demonstrate the proposed
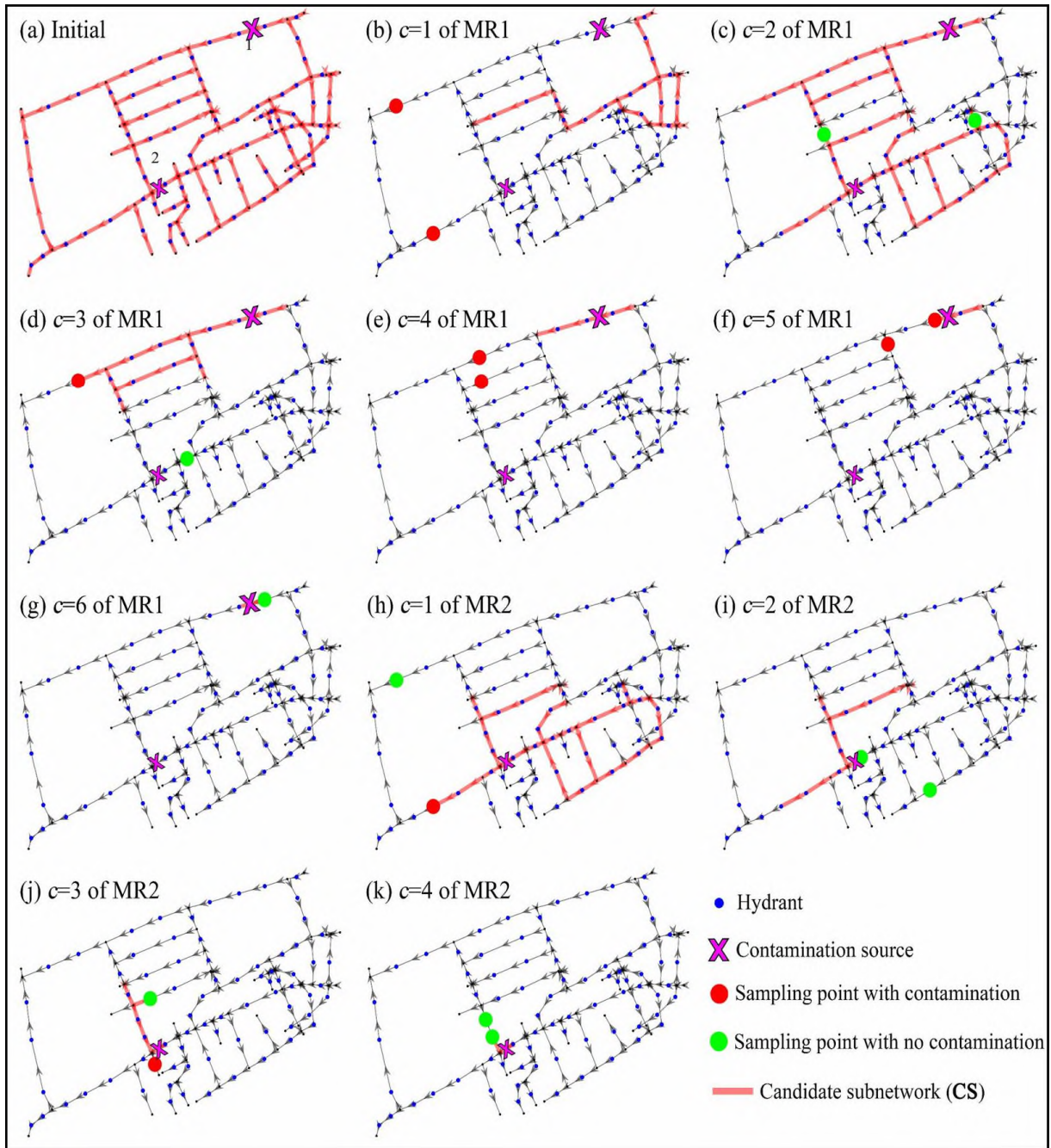437    method using two case studies at a DMA scale level.

438    For both case studies, we have analyzed a series of different combinations of sampling locations (i.e.,
439    the number of hydrants that can be simultaneously sampled) at each cycle, with $n$ ranging from 2
440    to 10. The number of potential contamination sources varies from one to three for these two WDSs.
441    The size of the MC method is determined to be 10,000 based on a preliminary analysis for both case
442    studies, but a larger value may be required for larger WDSs. The proposed MGSM is coded in C++
443    computing language with the aid of EPANET2.0 as the hydraulic solver to identify pipe flow
444    directions (He et al., 2018). For the DMA case study with 78 hydrants and two contamination sources,
445    the proposed method was tested using two and 10 potential sampling locations at each cycle required
446    an average of 102 and 54 seconds, respectively, on a PC with Intel i5-9400F CPU@2.90GHz. For
447    the MOD network with 143 hydrants and two contamination sources, the proposed MGSM with two
448    and 10 sampling locations at each cycle needs an average of 212 and 92 seconds, respectively. This
449    implies that the proposed method is very efficient to identify the optimal sampling locations based
450    on the test results. To enable the statistically rigorous analysis, for the single contamination source,
451    we considered all possible scenarios with one source assigned to each pipe of the network. For two
452    and three contamination sources, a total of 100 different randomly generated scenarios are
453    considered.

454    **4. Results and Discussion**

455     The proposed MGSM is demonstrated using the effectiveness (Section 4.1), the efficiency (Section

456     4.2) and the cost (Section 4.3) as shown in Section 4. The effectiveness is measured by the length of

457     finally identified pipes relative to the total pipe length of the entire WDS. The efficiency is measured

458     by the total number of sampling cycles, and the cost associated with the sampling process is

459     measured by the total number of samples that need to be tested in laboratory.

460     **4.1 Effectiveness of the proposed MGSM**

461     Figure 6 illustrates the application procedures of the proposed MGSM in dealing with the DMA case

462     study with two contamination sources (1 and 2 in Figure 6a) and two sampling locations at each

463     cycle. Two different MGSM runs (MR1 and MR2) are performed for this scenario, where the second

464     run assumed that the contamination source identified in the first run was eliminated. As shown in

465     this figure, in the beginning, the entire DMA is considered as the candidate sub-network (**CS**, Figure

466     6(a)) assuming that the water sample test at the outlet of this DMA shows contamination. This is

467     followed by the application of the MGSM, where six and four cycles were carried out to localize

468     contamination sources 1 and 2, respectively. The final identified pipe lengths associated with

469     contamination sources 1 and 2 are 741 and 762 meters, which represent only 1.26% and 1.30% of

470     the entire DMA, respectively. This implies that the proposed MGSM is able to effectively narrow

471     down the spatial range of pipes that contain contamination sources, which can greatly facilitate the

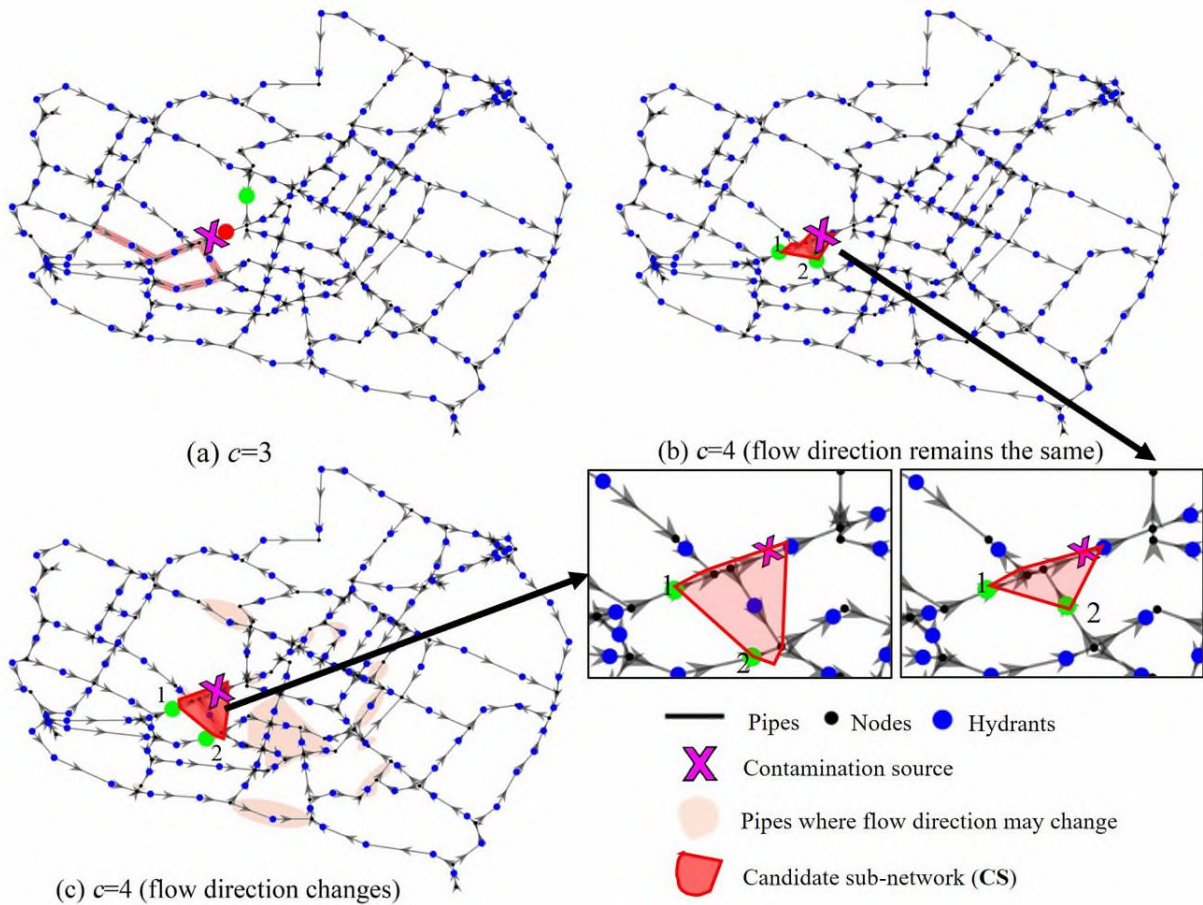472     subsequent field investigations to eliminate the cause of the problem.

**Figure 6: Source localization for the DMA case study with two contamination sources and two sampling locations at each cycle, where arrows indicate flow directions**
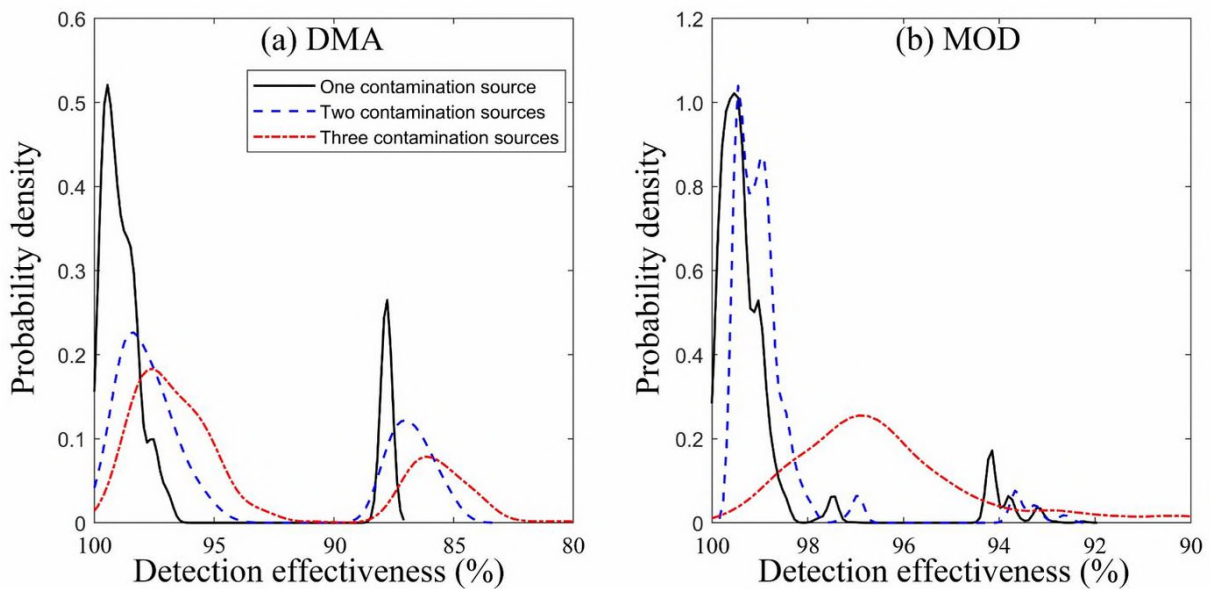
(a) *c*=3

(b) *c*=4 (flow direction remains the same)

**Pipes** ● **Nodes** ● **Hydrants**

✖ Contamination source

Pipes where flow direction may change

Candidate sub-network (**CS**)

(c) *c*=4 (flow direction changes)

**Figure 7: Source localization for the MOD case study with one contamination source and two sampling locations at each cycle, where arrows indicate flow directions**

Figure 7 illustrates the proposed MGSM applied to the WDS with possible pipe flow changes. As shown in this figure, if the pipe flow directions do not change, the two sampling locations identified by the proposed MGSM are 1 and 2 (Figure 7b) based on the candidate sub-network (**CS**) determined at $c$=3 (Figure 7a). However, if the flow directions change after the sample tests at $c$=3, the CS for the next cycle needs to account for such variation. For the given example, one pipe is added to the **CS** due to its flow changing. This addition affects the optimal sampling locations selected by the MGSM (the location of 2 is changed as shown in Figure 7c). Based on this example, the flow direction changes can be easily handled by the proposed MGSM. For the MOD case study, we assume the change in the flow direction status occurs (Figure 7c) after $c$=3, followed by a change to the original direction of flow after another two cycles.

23

489 It is found that the proposed MGSM is able to identify the contamination sources for all scenarios

490 considered in both case studies, implying its great effectiveness to localize contamination sources.

491 In this study, we define a detection effectiveness (%) metric as follows,

$$\text{Detection effectiveness} = (1 - \frac{L_f}{L_{all}}) \times 100\% \tag{6}$$

492 Where $L_f$ is the pipe length of the finally identified sub-network with contamination source(s) and

493 $L_{all}$ is the total pipe length of the entire WDS being considered. A high detection effectiveness

494 represents that the proposed method can greatly reduce the efforts or budgets of the subsequent field

495 investigations that are needed to micro-locate and eliminate contamination sources.



496

497 **Figure 8: Detection effectiveness (%) of the proposed MGSM applied to the two case studies**

498 Figure 8 presents the probability density of the detection effectiveness (%) for all contamination

499 scenarios considered, where the distribution of the ratio between the length of the finally identified

500 pipes and the total pipe length of the WDS for all contamination events is presented. It is seen from

501 this figure that the majority of the detection effectiveness (%) is higher than 95% and 98% for the

502 DMA and MOD case study respectively. This indicates that the finally identified pipes with

503 contamination source(s) represent a very small portion of the entire network, which can greatly

504 improve the efficiency of the subsequent engineering effort to fix the contamination problem. The

505 detection effectiveness (%) ranges between 80% and 90% for some contamination scenarios for the

24

506 DMA case study as shown in Figure 8(a). This is due to the sparse distribution of hydrants for these
507 events, and hence the length of the candidate sub-network identified by the proposed MGSM is
508 relatively large. The detection effectiveness (%) decreases when dealing with a larger number of
509 contamination sources that simultaneously exist in the WDS. It is noted that the detection
510 effectiveness (%) values are the same with those obtained using the average pipe length distance
511 between hydrants divided by the total pipe length of the network. This implies that the proposed
512 method is able to identify the pipe with contamination source between the two hydrants for each
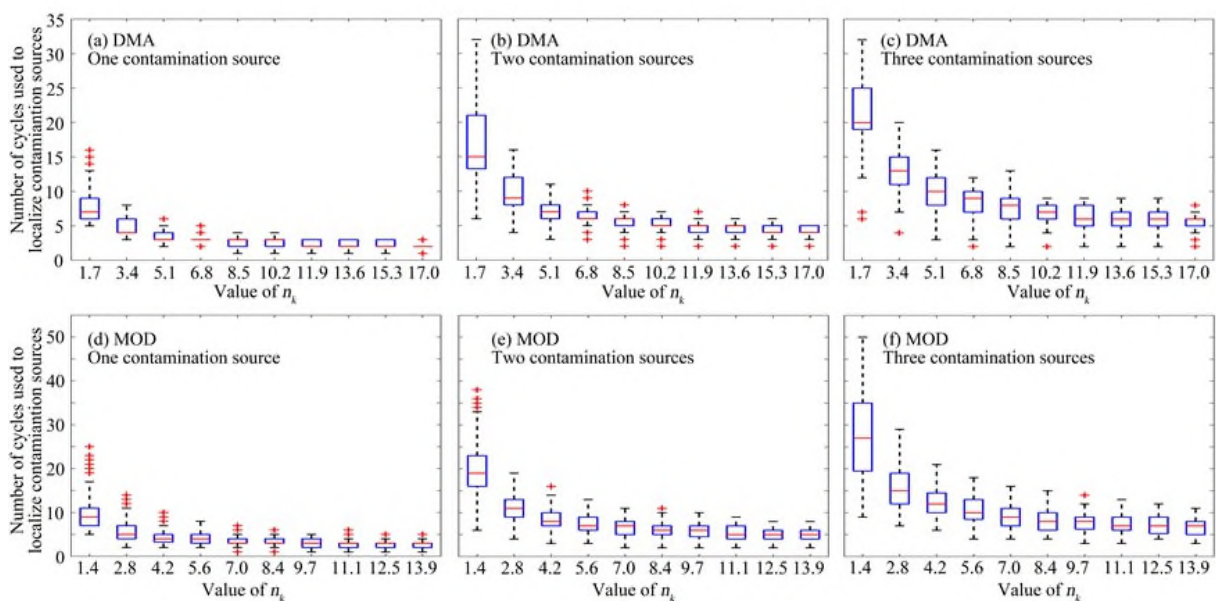513 scenario considered.

514 **4.2 Detection efficiency of the proposed MGSM**

515 The detection efficiency of the proposed MGSM can be evaluated using the number of total cycles
516 required for the entire procedure. The total time used in each cycle includes the time required to
517 collect and test samples, as well as the computation time needed to identify the sampling locations.
518 As previously stated, both the computation time and for sample collection are negligible compared
519 to the laboratory tests. Figure 9 shows the total number of cycles used to localize contamination
520 sources of the two case studies as a function of the varying number of samples per 100 km of pipe
521 length at each cycle ( $n_k$ ), where $n_k = n / L_{all} \times 100$. Such a normalization is used to enable the
522 generalization of the results to other WDSs.

523 As shown in Figure 9, an obvious trend that can be observed is that the detection efficiency is
524 improved when $n$ increases for all different contamination scenarios ( $n_k$ ranges from about 1.5 to 5)
525 for both case studies. A significant increase in efficiency occurs for $n_k$ >1.5, with improvements
526 diminishing when $n_k$>6. This is expected as a high $n_k$ value indicates a larger number of available
527 teams for collecting samples and a significant laboratory capacity for simultaneously testing multiple
528 samples. The diminishing efficiency improvement for large $n_k$ implies that an optimal sampling size
529 exists for the WDS when the efficiency is considered. For the DMA and MOD case studies, the
530 optimal $n_k$ value can be between 7.0 and 8.5 as a further increase in $n_k$ value does not significantly
531 improve the MGSM's detection efficiency, as shown in Figure 9. However, the optimal $n_k$ value for
532 detection efficiency can be case study dependent as it can be related to the size of the WDS being

25

533 considered. In addition, a large $n_k$ value corresponds to a significant financial commitment, and

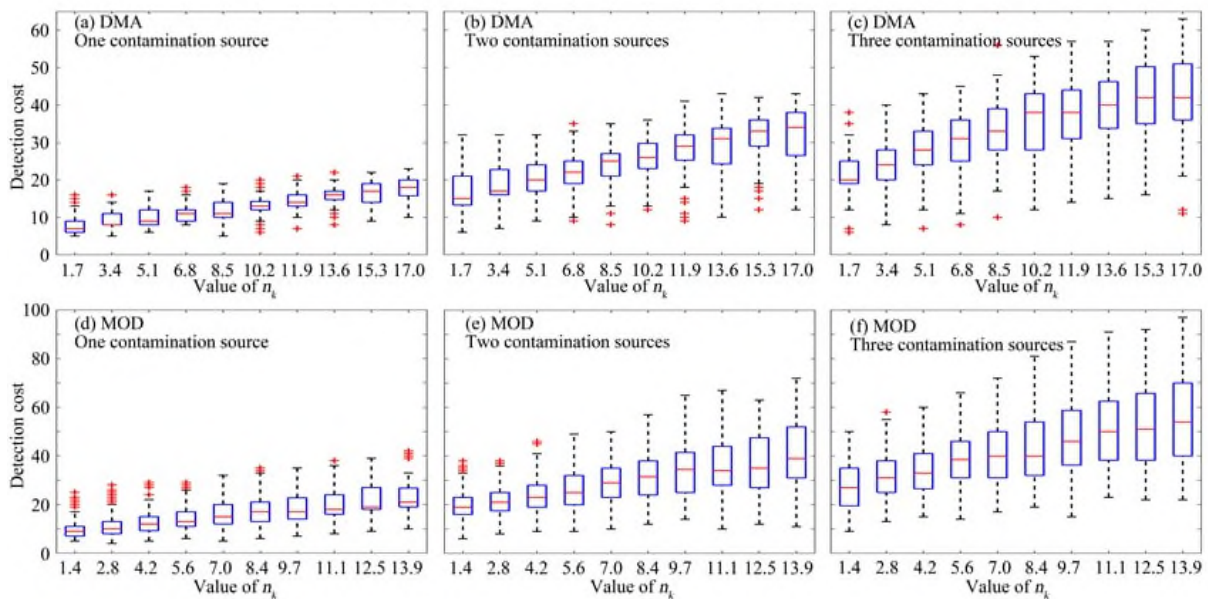534 hence the decision process can be also affected by the budgets available.

535 Interestingly, for the same number of sampling locations at each 100 km pipe length $n_k$ , when $n_k$ is

536 relatively low, the total number of cycles can vary significantly. For example, for the DMA case

537 study if $n_k$=1.7, the detection efficiency can vary from 5 to 15 cycles for the one contamination

538 source, and range from 7 to 25 cycles when three contamination sources are simultaneously

539 considered. Similar observations can be made for the MOD case study. This implies that the location

540 of the contamination sources can appreciably affect the detection efficiency when there is a low

541 number of sampling teams available and/or a limited laboratory capacity for testing multiple samples.

542 When a sufficiently large $n_k$ is considered, the detection efficiency variations become small, as

543 observed in Figure 9. This implies that the choice of $n_k$ will also affect the uncertainty associated

544 with method efficiency, which should be also accounted for in engineering practice.



545

546 **Figure 9: The number of cycles used to localize contamination sources versus the number of**

547 **sampling points for every 100 km pipe length at each cycle ($n_k$) for the proposed MGSM**

548 **applied to the two case studies**

549 **4.3 Detection cost of the proposed MGSM**

550 In this study, the detection cost of the proposed MGSM is measured by the total number of samples
551 that have been tested to localize the contamination sources. Figure 10 shows the detection cost as a
552 function of varying $n_k$ for both case studies. Despite some variations, a large $n_k$ value is generally
553 associated with a greater detection cost for both case studies. In addition, the simultaneous presence
554 of a larger number of contamination sources also causes an overall increase in detection costs. This
555 information combined with the efficiency results in Figure 9 can be used as guidance for developing
556 effective water quality sampling plans or budgets for a given WDS.



557

558 **Figure 10: Detection cost (i.e., the number of total samples) versus the number of sampling**
559 **points for every 100 km pipe length at each cycle ( $n_k$ ) for the proposed MGSM applied to the**
560 **two case studies**

561 **5. Summary and Conclusions**

562 Existing research on water quality management and contamination source localization in WDSs has
563 focused mainly on developing methods that assume availability of accurate water quality models
564 and multi-parameter online sensors. However, that is not true for many water utilities. A promising
565 way to address such problems is through the iterative manual grab-sample strategies, thereby
566 enabling effective contaminant localizing. To this end, this study proposes a new method for water

27

567  quality manual grab-sampling (termed as MGSM in this paper) to enable identification of
568  contamination sources in WDSs.

569  The proposed MGSM is suitable for situations where online multi-parameter water quality sensors
570  are sparsely available or completely missing, which is the case with many utilities. This is mainly
571  due to the high purchase and maintenance cost associated with these sensors, as well as their inability
572  (or inaccurate) to detect the complex water quality parameters (e.g., metals, microorganism and
573  personal care products, Jia et al. (2021b)). In addition, a grab-sampling method is tailored for the
574  cases when contamination is *continually* present in the WDS and with slow or low impacts to the
575  WDSs. That is the case with misconnections between water supply pipes and sewer (or grey) pipes
576  and contaminations caused by pipe leaks, corrosion or hydraulic turbulence. For events with serious
577  consequences, the candidate sub-networks (CSs) with contamination sources may need to be shut
578  down or sampled manually as much as possible.

579  Based on the results obtained for two real-world cases, the following findings and conclusions can
580  be drawn:
581  (1)    The newly proposed MGSM can successfully detect and locate continuous contamination
582      source(s) for a wide range of scenarios, including multiple contamination source(s) in complex
583      WDSs with varying pipe flow directions. This is a significant advantage over the traditional
584      approach that works only with one contamination source and fixed flow directions, as described
585      in Wong et al. (2010).
586  (2)    For the two case studies, the new MGSM identified contamination source(s) within 5% of
587      the total pipe length of the WDS. This indicates the high effectiveness of the proposed MGSM in
588      narrowing narrow down the spatial range of the sub-network with potential contamination sources.
589      From the practical point of view, it also improves the efficiency of maintenance efforts to eliminate
590      the sources of contamination.
591  (3)    The detection efficiency (measured by the number of sampling and testing cycles) of the
592      MGSM can be significantly improved when the number of sampling points per 100 km pipe length
593      at each cycle ( $n_k$ ) increases from about 1.5 to a moderate value (e.g. $n_k \approx 7$). The increase in
594      efficiency diminishes with further increases in $n_k$. This implies that there exists an optimal $n_k$ value
595      for a given WDS, representing the balanced trade-off between detection efficiency and costs

596    associated with methodology. The detection cost grows with the increase in the number of sampling

597    points per 100 pipe length, $n_k$. All these findings are important for the implementation of the

598    method as they can guide the process of selecting the optimal number of sampling teams and

599    required laboratory capacity.

600    In view of the practical application, the proposed MGSM can be used to regularly check water

601    quality safety for WDSs with a low density of sensors as this is routine work in many water utilities.

602    For instance, in China, many water utilities need to take water samples from hydrants or end users

603    every month, with the number of samples depending on the scale of the WDS and importance level

604    of the city. These water samples are comprehensively measured in the laboratory following the

605    Water Quality Standard that has 106 parameters. Many water utilities collect grab samples from

606    large WDSs at fixed locations based on specialists' engineering expertise. For example, a

607    practitioner may collect grab samples from all established fixed locations (if say, 50 locations) and

608    test for a combination (or all) of the specified water quality parameters in the laboratory. Such a

609    strategy is time-consuming and expensive (labor and measurement costs). Therefore, the sampling

610    strategy can be improved with the aid of the proposed MGSM in order to save the cost. It can be

611    concluded that the MGSM is an alternative to the sensor-based detection methods.

612    The limitation of the method proposed here is the potentially high cost and time required to identify

613    the source(s) as all grab samples need to be collected manually (with technicians moving between

614    different locations during multiple cycles) and processed (in the lab). In addition, the pipes identified

615    as the potential contamination sources need to be inspected in the field to micro-locate the

616    contamination source(s) with the aid of manual checking or detection robots (Huang et al., 2020).

617    This too requires time and has a cost associated with it. This, however, applies to most of the existing

618    sensor-based methods as well. Another limitation is that the proposed MGSM can be only applicable

619    to contamination events with continuous injections to the WDS conditioned on known pipe flow

620    directions. Furthermore, when dealing with scenarios with pipe flow changes, there is likely that

621    such changes would affect the utility of the proposed MGSM, which needs attention within practical

622    implementation. While the practical application of the developed MGSM can be simple as it only

623    requires flow direction information (Zhang et al., 2021), it should be also acknowledged the flow

624    information can be challenging for some old pipes due to system uncertainties.

Future studies along this research line include (i) the application of the proposed method to further large real WDSs; (ii) the extension of the graph partitioning strategy within the proposed MGSM to account for both the pipe length and pipe velocity; (iii) the extension of the proposed MGSM to deal with contamination events with intermittent injections to the WDS.

**Acknowledgements**

**References**

Asheri Arnon, T., Ezra, S. and Fishbain, B. (2019). Water characterization and early contamination detection in highly varying stochastic background water, based on Machine Learning methodology for processing real-time UV-Spectrophotometry. *Water Res.* 155, 333-342.

Bhatia, R. and Davis, C. (1995). A Cauchy-Schwarz inequality for operators with applications. Linear Algebra Appl. 223-224, 119-129.Butera, I., Gómez-Hernández, J.J. and Nicotra, S. 2021. Contaminant-Source Detection in a Water Distribution System Using the Ensemble Kalman Filter. *J. Water Resour. Plan. Manag.* 147(7), 04021029.

ChinaNews, (2020). http://www.chinanews.com/sh/2020/07-30/9252169.shtml.

Di Nardo, A., Giudicianni, C., Greco, R., Herrera, M., Santonastaso, G.F. and Scala, A. (2018). Sensor placement in water distribution networks based on spectral algorithms. *13th International Conference on Hydroinformatics (HIC2018)* 7.

Giudicianni, C., Herrera, M., Nardo, A.D., Greco, R., Creaco, E. and Scala, A. (2020). Topological Placement of Quality Sensors in Water-Distribution Networks without the Recourse to Hydraulic Modeling. *J. Water Resour. Plan. Manag.* 146(6), 04020030.

Grbčić, L., Lučin, I., Kranjčević, L. and Družeta, S. (2020). Water supply network pollution source identification by random forest algorithm. *J. Hydroinformatics* 22(6), 1521-1535.

Hart, D., Rodriguez, J.S., Burkhardt, J., Borchers, B., Laird, C., Murray, R., Klise, K. and Haxton, T. (2019). Quantifying Hydraulic and Water Quality Uncertainty to Inform Sampling of Drinking Water Distribution Systems. *J. Water Resour. Plan. Manag.* 145(1), 04018084.

Hart, W.E. and Murray, R. (2010). Review of Sensor Placement Strategies for Contamination Warning Systems in Drinking Water Distribution Systems. *J. Water Resour. Plan. Manag.* 136(6), 611-619.

He, G., Zhang, T., Zheng, F. and Zhang, Q. (2018). An efficient multi-objective optimization method for water quality sensor placement within water distribution systems considering contamination probability variations. *Water Res.* 143, 165-175.

He, G., Zhang, T., Zheng, F., Li, C., Zhang, Q., Dong, F. and Huang, Y. (2019) Reaction of fleroxacin with chlorine and chlorine dioxide in drinking water distribution systems: Kinetics, transformation mechanisms and toxicity evaluations. *Chem. Eng. J.* 374, 1191-1203.

Hu, C., Dai, L., Yan, X., Gong, W., Liu, X. and Wang, L. (2020). Modified NSGA-III for sensor placement in water distribution system. *Inf. Sci.* 509, 488-500.

Hu, C., Ren, G., Liu, C., Li, M. and Jie, W. (2017). A Spark-based genetic algorithm for sensor placement in large scale drinking water distribution systems. *Cluster Comput.* 20(2), 1089-1099.

Hu, C., Zhao, J., Yan, X., Zeng, D. and Guo, S. (2015). A MapReduce based Parallel Niche Genetic Algorithm for contaminant source identification in water distribution network. *Ad Hoc Netw.* 35, 116-126.

Huang, Y., Zheng, F., Kapelan, Z., Savic, D., Duan, H.-F. and Zhang, Q. (2020) Efficient Leak Localization in Water Distribution Systems Using Multistage Optimal Valve Operations and Smart Demand Metering. *Water Resour. Res.* 56(10), e2020WR028285.

Jerez, D.J., Jensen, H.A., Beer, M. and Broggi, M. (2021). Contaminant source identification in water distribution networks: A Bayesian framework. *Mech. Syst. Signal Process* 159, 107834.

Jia, Y., Zheng, F., Zhang, Q., Duan, H.-F., Savic, D. and Kapelan, Z. (2021a). Foul sewer model development using geotagged information and smart water meter data. *Water Res.* 204, 117594.

Jia, Y., Zheng, F., Maier, H.R., Ostfeld, A., Creaco, E., Savic, D., Langeveld, J. and Kapelan, Z. (2021b) Water quality modelling in sewer networks: review and future research directions. Water Research, 117419.

685 Khorshidi, M.S., Nikoo, M.R. and Sadegh, M. (2018). Optimal and objective placement of sensors
686   in water distribution systems using information theory. *Water Res.* 143, 218-228.

687 Li, C., Yang, R., Zhou, L., Zeng, S., Mavrovouniotis, M., Yang, M., Yang, S. and Wu, M. (2021).
688   Adaptive Multipopulation Evolutionary Algorithm for Contamination Source Identification
689   in Water Distribution Systems. *J. Water Resour. Plan. Manag.* 147(5), 04021014.

690 Maier, H.R., Kapelan, Z., Kasprzyk, J., Kollat, J., Matott, L.S., Cunha, M.C., Dandy, G.C., Gibbs,
691   M.S., Keedwell, E., Marchi, A., Ostfeld, A., Savic, D., Solomatine, D.P., Vrugt, J.A., Zecchin,
692   A.C., Minsker, B.S., Barbour, E.J., Kuczera, G., Pasha, F., Castelletti, A., Giuliani, M. and
693   Reed, P.M. (2014) Evolutionary algorithms and other metaheuristics in water resources:
694   Current status, research challenges and future directions. *Environ. Model Softw.* 62(0), 271-
695   299.

696 Mann, A.V., McKenna, S.A., Hart, W.E. and Laird, C.D. (2012). Real-time inversion in large-
697   scale water networks using discrete measurements. *Comput. Chem. Eng.* 37, 143-151.

698 Naserizade, S.S., Nikoo, M.R. and Montaseri, H. (2018). A risk-based multi-objective model for
699   optimal placement of sensors in water distribution system. *J. Hydrol.* 557, 147-159.

700 Ohar, Z., Lahav, O. and Ostfeld, A. (2015). Optimal sensor placement for detecting
701   organophosphate intrusions into water distribution systems. *Water Res.* 73, 193-203.

702 Oliker, N. and Ostfeld, A. (2014). A coupled classification – Evolutionary optimization model for
703   contamination event detection in water distribution systems. *Water Res.* 51, 234-245.

704 Ostfeld, A., Uber, J.G., Salomons, E., Berry, J.W., Hart, W.E., Phillips, C.A., Watson, J.-P., Dorini,
705   G., Jonkergouw, P., Kapelan, Z., Pierro, F.d., Khu, S.-T., Savic, D., Eliades, D., Polycarpou,
706   M., Ghimire, S.R., Barkdoll, B.D., Gueli, R., Huang, J.J., McBean, E.A., James, W., Krause,
707   A., Leskovec, J., Isovitsch, S., Xu, J., Guestrin, C., VanBriesen, J., Small, M., Fischbeck, P.,
708   Preis, A., Propato, M., Piller, O., Trachtman, G.B., Wu, Z.Y. and Walski, T. 2008. The Battle
709   of the Water Sensor Networks (BWSN): A Design Challenge for Engineers and Algorithms.
710   J. Water Resour. Plan. Manag. 134(6), 556-568.

711 Ostfeld, A., Oliker, N. and Salomons, E. (2014). Multiobjective Optimization for Least Cost
712   Design and Resiliency of Water Distribution Systems. *J. Water Resour. Plan. Manag.*
713   140(12), 04014037.

714 Preis, A. and Ostfeld, A. (2006). Contamination Source Identification in Water Systems: A Hybrid
715   Model Trees–Linear Programming Scheme. *J. Water Resour. Plan. Manag.* 132(4), 263-273.

716 Preis, A. and Ostfeld, A. (2007). A contamination source identification model for water
717     distribution system security. *Eng. Optim.* 39(8), 941-947.

718 Preis, A. and Ostfeld, A. (2008). Genetic algorithm for contaminant source characterization using
719     imperfect sensors. *Civ. Eng. Environ. Syst.* 25(1), 29-39.

720 Qi, Z., Zheng, F., Guo, D., Maier, H.R., Zhang, T., Yu, T. and Shao, Y. (2018). Better
721     Understanding of the Capacity of Pressure Sensor Systems to Detect Pipe Burst within Water
722     Distribution Networks. *J. Water Resour. Plan. Manag.* 144(7), 04018035.

723 Rathi, S. and Gupta, R. (2014). Sensor Placement Methods for Contamination Detection in Water
724     Distribution Networks: A Review. *Procedia Eng.* 89, 181-188.

725 Robertson, L., Gjerde, B., Hansen, E.F. and Stachurska-Hagen, T. (2008). A water contamination
726     incident in Oslo, Norway during October 2007; a basis for discussion of boil-water notices
727     and the potential for post-treatment contamination of drinking water supplies. *J. Water Health*
728     7(1), 55-66.

729 Rodriguez, J.S., Bynum, M., Laird, C., Hart, D.B., Klise, K.A., Burkhardt, J. and Haxton, T.,
730     (2021). Optimal Sampling Locations to Reduce Uncertainty in Contamination Extent in
731     Water Distribution Systems. *J. Infrastruct. Syst.* 27(3), p.04021026.

732 Sankary, N. and Ostfeld, A. (2018). Multiobjective Optimization of Inline Mobile and Fixed
733     Wireless Sensor Networks under Conditions of Demand Uncertainty. *J. Water Resour. Plan.*
734     *Manag.* 144(8), 04018043.

735 Sankary, N. and Ostfeld, A. (2019). Bayesian Localization of Water Distribution System
736     Contamination Intrusion Events Using Inline Mobile Sensor Data. *J. Water Resour. Plan.*
737     *Manag.* 145(8), 04019029.

738 Sun, L., Yan, H., Xin, K. and Tao, T. (2019). Contamination source identification in water
739     distribution networks using convolutional neural network. *Environ. Sci. Pollut. Res.* 26(36),
740     36786-36797.

741 Tinelli, S., Creaco, E. and Ciaponi, C. (2017). Sampling Significant Contamination Events for
742     Optimal Sensor Placement in Water Distribution Systems. *J. Water Resour. Plan. Manag.*
743     143(9), 04017058.

744 Ung, H., Piller, O., Gilbert, D. and Mortazavi, I. (2017). Accurate and Optimal Sensor Placement
745     for Source Identification of Water Distribution Networks. *J. Water Resour. Plan. Manag.*
746     143(8), 04017032.

747 Vrachimis, S.G., Lifshitz, R., Eliades, D.G., Polycarpou, M.M. and Ostfeld, A. (2020). Active
748     Contamination Detection in Water-Distribution Systems. *J. Water Resour. Plan. Manag.*
749     146(4), 04020014.

750 Weickgenannt, M., Kapelan, Z., Blokker, M. and Savic, D.A. (2010). Risk-Based Sensor
751     Placement for Contaminant Detection in Water Distribution Systems. *J. Water Resour. Plan.*
752     *Manag.* 136(6), 629-636.

753 Winter, C.d., Palleti, V.R., Worm, D. and Kooij, R. (2019). Optimal placement of imperfect water
754     quality sensors in water distribution networks. *Comput. Chem. Eng.* 121, 200-211.

755 Wong, A., Young, J. and Laird, C.D. (2010). Optimal determination of grab sample location and
756     source inversion in large-scale water distribution systems. *Water Distribution Systems*
757     *Analysis 2010.*

758 Yang, X. and Boccelli, D.L. (2014). Bayesian Approach for Real-Time Probabilistic
759     Contamination Source Identification. *J. Water Resour. Plan. Manag.* 140(8), 04014019.

760 Yang, X. and Boccelli, D.L. (2016). Model-Based Event Detection for Contaminant Warning
761     Systems. *J. Water Resour. Plan. Manag.* 142(11), 04016048.

762 Zhang, Q., Zheng, F., Jia, Y., Savic, D. and Kapelan, Z. (2021) Real-time foul sewer hydraulic
763     modelling driven by water consumption data from water distribution systems. *Water Res.* 188,
764     116544.

765 Zhang, Q., Zheng, F., Kapelan, Z., Savic, D., He, G. and Ma, Y. (2020) Assessing the global
766     resilience of water quality sensor placement strategies within water distribution systems.
767     *Water Res.* 172, 115527.

768 Zheng, F., Du, J., Diao, K., Zhang, T., Yu, T. and Shao, Y. (2018). Investigating Effectiveness of
769     Sensor Placement Strategies in Contamination Detection within Water Distribution Systems.
770     *J. Water Resour. Plan. Manag.* 144(4), 06018003.