

1 **Title: Machine Learning to Access and Ensure Safe Drinking Water Supply: A Systematic**

2 **Review**

3 Feng Feng^{1,2†}, Zhenru Chen^{3†}, Jianyuan Ni⁴, Yuanxun Zhang¹, Yuan Feng⁵, Yunchao Xie^{6*}, Chiqian
4 Zhang^{7*}

5
6 **Affiliations:** ¹Department of Electrical Engineering and Computer Science, University of Missouri,
7 Columbia, Missouri 65211, United States

8 ²Department of Developmental Neurobiology, St. Jude Children’s Research Hospital, Memphis,
9 Tennessee 38105, United States

10 ³Department of Mechanical and Aerospace Engineering, University of Missouri, Columbia,
11 Missouri 65211, United States

12 ⁴Department of Computer Science, Lamar University, Beaumont, Texas 77710, United States

13 ⁵Department of Computational Biology, St. Jude Children’s Research Hospital, Memphis,
14 Tennessee 38105, United States

15 ⁶Department of Mechanical and Manufacturing Engineering, Miami University, Oxford, Ohio
16 45056, United States

17 ⁷Civil Engineering Program, College of Engineering & Computer Science, Arkansas State
18 University, Arkansas 72467, United States

19

20 †These authors contributed equally

21 *E-mails: czhang@astate.edu; xiey54@miamioh.edu

22 **Abstract** Drinking water is essential to public health and socioeconomic growth. Therefore,
23 assessing and ensuring drinking water supply is a critical task in modern society. Conventional
24 approaches to analyzing and controlling drinking water quality are labor-intensive and costly with
25 a low throughput. Machine learning (ML) is an alternative, promising technique for assessing and
26 ensuring safe drinking water supply. Existing reviews have summarized the applications of ML in
27 safe drinking water supply from different aspects. However, a state-of-the-art, comprehensive
28 review is missing that focuses on applying ML to monitor, simulate, predict, and control drinking
29 water quality, especially in engineered water systems. This review, therefore, critically compiles
30 the applications of ML in assessing and ensuring water quality in engineered water systems. To be
31 comprehensive, we also cover the applications of ML in other drinking-water-related settings such
32 as water sources and water purification processes. We explain the basic mechanics and workflows
33 of ML and focus on the applications of ML to access and control key factors or etiologies in
34 drinking water from the physical, chemical, and microbiological aspects. Those factors or
35 etiologies affect water quality and public health, such as water pipeline failures, disinfectant by-
36 products, heavy metals, opportunistic pathogens, biofilms, and antimicrobial resistance genes. We
37 then present a macroscopic illustration to display the distribution of ML models across research
38 topics in safe drinking water supply. Neural-network-based and regression-based models are the
39 top two models frequently used in the field of drinking water supply. We finally discuss the
40 challenges and outlooks for the applications of machine learning in safe drinking water supply.
41 Filling the gap between the water research and the AI research communities and using AI to solve
42 the global drinking water crisis should be the main focus of future research. This is the first review
43 summarizing the feasibility and applications of ML in assessing and ensuring water quality in
44 municipal engineered water systems as well as related water environments.

45 **Keywords:** Drinking water quality; Engineered water systems; Artificial intelligence;
46 Opportunistic pathogens; Disinfection byproducts; Heavy metals
47

48 1. Introduction

49 Clean and safe drinking water is vital to public health(Cabral, 2010; Cutler and Miller, 2005; Li
50 and Wu, 2019). Most people obtain drinking water from public water systems. For instance, over
51 90% of people in the United States (U.S.) access drinking water from approximately 150,000
52 public water systems(U.S. EPA, 2022). Public water systems provide water for human
53 consumption through engineered water systems (including drinking water distribution systems or
54 DWDSs and building premise plumbing systems) (Zhang and Lu, 2021b). Drinking-water-related
55 disease outbreaks and chronic diseases cause many hospitalizations and deaths, leading to
56 significant socioeconomic losses (Benedict et al., 2017; Craun et al., 2010; Lee et al., 2023).
57 Therefore, assessing and ensuring water quality, especially in engineered water systems, is critical
58 to the health of the end consumers (WHO, 2011).

59
60 Assessing and ensuring water quality in engineered water systems are complex(Li and Wu, 2019).
61 Multiple ever-changing variables affect water quality in engineered water systems such as the
62 quality of water sources, treatment processes, water pipe materials, system configuration and
63 length, natural disasters, and geographical factors (Delpla et al., 2009; Li and Wu, 2019; Proctor
64 et al., 2020). For instance, the effluent at water utilities many have high quality, while the quality
65 deteriorates in engineered water systems because of microbial (re)growth, the formation of
66 disinfection by-products (DBPs), pipe failures (e.g., breaks and leaks), and the detachment of
67 heavy metals from pipes(Li et al., 2019; Liu et al., 2016). Meanwhile, engineered water systems
68 have various hazardous agents such as harmful microbes (especially opportunistic pathogens or
69 OPs) (Zhang et al., 2021), DBPs(Benítez et al., 2021; Lee et al., 2013), heavy metals(Chowdhury
70 et al., 2016; Gonzalez et al., 2013), pesticides and herbicides(Syafudin et al., 2021)

71 (Mukhopadhyay et al., 2022), and other emerging contaminants (e.g., antimicrobials and
72 microplastics) (Gogoi et al., 2018; Kirstein et al., 2021; Taheran et al., 2018). Those agents are
73 interconnected, and controlling only one group of agents frequently fails to secure the drinking
74 water quality. For instance, increasing disinfectant residual concentrations suppresses microbial
75 regrowth in engineered water systems but promotes the formation of DBPs(Zhang and Lu, 2021c).
76 By contrast, reducing the dose of disinfectant residuals in engineered water systems can mitigate
77 the DBP issue, but microbes (including pathogens) can thrive.

78
79 Because of the complex nature of water quality in engineered water systems, assessing and
80 ensuring drinking water quality using conventional means is challenging. Traditional methods are
81 time-consuming, labor-intensive, inefficient (i.e., low throughput), and costly(Ahmed et al., 2019)
82 (Zainurin et al., 2022). Artificial intelligence (AI), especially machine learning (ML), is promising
83 to address the deficiencies in the traditional approaches to access and ensure safe drinking water
84 supply (Richards et al., 2023). The adaptability and predictive power of ML offer significant
85 advantages over other AI technologies (Willard et al., 2022; Zhu et al., 2022), particularly when
86 handling drinking water quality with a dynamic and complex nature. Consequently, ML emerges
87 as a specialized application of AI to enhance drinking water treatment and quality(Henrique Alves
88 Ribeiro and Reynoso-Meza, 2023; Li et al., 2021; Narita et al., 2023; Speight et al., 2019).

89
90 Existing reviews have summarized the applications of ML in various aspects of drinking water
91 quality (Ewuzie et al., 2022; Huang et al., 2021; Zhu et al., 2022), such as source water quality and
92 contamination(Gong et al., 2023; Zanoni et al., 2022), drinking water treatment(Li et al., 2021;
93 Lowe et al., 2022; Ortiz-Lopez et al., 2022), and drinking water quality anomaly detection (Dogo

94 et al., 2019). However, a state-of-the-art, comprehensive review is missing that focuses on
95 applying ML to monitor, simulate, predict, and control drinking water quality, especially in
96 engineered water systems. Since engineered water systems are the vital civil infrastructure
97 delivering municipal water from water utilities to the residents and industrial/commercial
98 consumers (WHO, 2011), summarizing such applications can help ensure drinking water quality,
99 protect public health, and promote socioeconomic development.

100

101 In this review, we critically compile the applications of ML in assessing and ensuring water quality
102 across all stages of drinking water treatment and distribution (especially in engineered water
103 systems) from the physical, chemical, and microbiological aspects (Figure 1). We focus on the
104 applications of ML to access and control key factors or etiologies in drinking water that affect
105 water quality and public health, such as water pipeline failures, DBPs, heavy metals, OPs, biofilms,
106 and antimicrobial resistance genes (ARGs). Since water sources and the treatment processes
107 govern drinking water quality, we also discuss the applications of ML in monitoring source water
108 quality and the efficiency of water treatment technologies. This is the first review focusing on the
109 applications of ML in assessing and controlling water quality in engineered water systems, while
110 the applications of ML in other drinking-water-related settings are also summarized.

111

112 **2. Literature search strategy and inclusion/exclusion criteria**

113 We follow previous protocols to search the literature on the applications of ML to assess and ensure
114 water quality in engineered water systems and related settings (Tolaymat et al., 2010) (Miake-Lye
115 et al., 2016; Mostafavifar et al., 2012; Zhang et al., 2016; Ziegelbauer et al., 2012). We mainly
116 search Google Scholar (scholar.google.com) but also include other databases such as

117 ScienceDirect (sciencedirect.com), PubMed Central[®] (ncbi.nlm.nih.gov/pmc/), and ACS
118 Publications (pubs.acs.org). We include only high-quality English articles published in or before
119 December 2023 in top-tier, peer-reviewed academic journals, books, and conference proceedings.
120 We manually review the retrieved publications and retain only those focusing on the applications
121 of ML in assessing and ensuring drinking water quality. We also manually check articles that cited
122 and were cited by those retrieved publications and keep the most relevant ones for further analysis.

123
124 The main literature search keywords are “machine learning,” “deep learning,” “artificial
125 intelligence,” “models,” “modeling,” “prediction,” “networks,” “anomaly detection,” “municipal,”
126 “water sources” “water treatment,” “water utilities,” “water plants,” “water treatment facilities,”
127 “distribution systems,” “building,” “premise plumbing,” “pipes,” “pipelines,” “water mains,”
128 “drinking water,” “potable water,” “tap water,” “municipal water,” “bulk water,” “disinfection,”
129 “disinfectants,” “free chlorine,” “chlorine,” “monochloramine,” “disinfectant residuals,”
130 “disinfection by-products,” “heavy metals,” “lead,” “copper,” “nitrification,” “pipe failures,”
131 “water-related,” “water-borne,” “disease outbreaks,” “public health,” “microbes,” “bacteria,”
132 “biofilms,” “opportunistic pathogens,” “opportunistic premise plumbing pathogens,” “*Legionella*,”
133 “*Mycobacterium*,” “antibiotic resistance,” “antimicrobial resistance,” “antimicrobial resistance
134 genes,” “*Cryptosporidium*,” and “*Giardia*.”

135

136 3. Machine learning primers

137 In the past decade, ML has sparked substantial progress across domains in modern society,
138 including object detection (Erhan et al., 2014; Lin et al., 2017), autonomous driving (Almalioglu
139 et al., 2022; Feng et al., 2023), drug delivery (Allesoe et al., 2023; Wołos et al., 2022), playing

140 games (Kaufmann et al., 2023; Vinyals et al., 2019), weather forecasting (Bi et al., 2023), and
141 design-by-analogy (Jiang et al., 2021). Three key advancements in AI and computer science drive
142 this process: I) the availability of extensive datasets, II) the development of robust computing
143 hardware, and III) the refinement of advanced algorithms. In contrast to traditional physical and
144 chemical theories relying on explicit formulas for problem-solving, ML tackles problems by
145 extracting concealed insights from datasets through the learning process (Ley et al., 2022).

146

147 **3.1 Categories of machine learning**

148 On the basis of the quantity and nature of available datasets, ML can be divided into four main
149 categories: I) supervised learning, II) unsupervised learning, III) semi-supervised learning, and IV)
150 reinforcement learning (Goodfellow et al., 2016). Supervised learning uses both input data and
151 corresponding labels for training. Unsupervised learning, on the other hand, deals solely with input
152 data without labeled information. As a combination of supervised and unsupervised learning, semi-
153 supervised learning combines mostly unlabeled datasets and a limited number of labeled ones.
154 Reinforcement learning algorithms, such as Q-learning, enable learning by interacting with an
155 environment and receiving feedback. These algorithms underpin the growth of AI, improving
156 system performance through exposure to data and experience. Among the four categories,
157 supervised learning is the most widely used and well-established in assessing and protecting
158 drinking water quality because of its strength in prediction with labeled datasets (Cordero et al.,
159 2021; Hong et al., 2020; Zhang et al., 2019; Zhou et al., 2019). Conversely, the application of semi-
160 supervised learning in safe drinking water supply is scarce, and reinforcement learning is
161 unexplored.

162 Supervised learning is acquiring a mapping function from the input data to the corresponding
163 output data on the basis of a labeled set of input-output pairs or conditional distributions
164 (Goodfellow et al., 2016). During training, the algorithm adjusts its parameters to minimize the
165 discrepancy between predicted and actual outputs. Supervised learning is widely used in tasks such
166 as classification and regression. Examples of supervised learning models include naïve Bayes (NB)
167 (Gomez-Alvarez and Revetta, 2020), logistic regression (LR) (Bagriacik et al., 2018), support
168 vector machines (SVM) (Oh et al., 2021), k-nearest neighbor (KNN) (Ghiassi et al., 2017),
169 decision trees (DT) (Shi et al., 2022), random forests (RF) (Berglund et al., 2023; Cha et al., 2021),
170 and extreme gradient boosting (XGB) (Park et al., 2020).

171
172 Unsupervised learning focuses on extracting patterns, structures, and relationships from the input
173 data without labeled outputs (Goodfellow et al., 2016). Instead of targeting predefined targets,
174 unsupervised learning algorithms (such as K-means and dimensionality reduction techniques)
175 discover inherent structures within the data. For instance, K-means groups similar data points on
176 the basis of their intrinsic features (Moodley and van der Haar, 2019). Dimensionality reduction
177 techniques including principal component analysis (PCA) simplify complex datasets by preserving
178 their essential characteristics (Peleato et al., 2018). In drinking water research, unsupervised
179 learning is crucial in tasks such as clustering bacteria (Moodley and van der Haar, 2019; Pinto et
180 al., 2014), simplifying data for subsequent analysis (Peleato et al., 2018), and analyzing raw data
181 to identify key parameters and major relationships affecting water quality (Kazemi et al., 2023).

182

183 **3.2 Workflow chart of machine learning**

184 **3.2.1 Define the problem**

185 Defining the problem is critical in applying ML that converts a complex challenge into a well-
186 defined scope and purpose. To start, one should define the objectives, outline desired outcomes,
187 and determine if the task is a regression problem (such as predicting DBP concentration) or a
188 classification problem (such as categorizing drinking water contamination status and pipe burst
189 localization).

190

191 **3.2.2 Gather data**

192 In any ML endeavor, the quality of data is the key to the success of the subsequent modeling. Data
193 collection involves sourcing, gathering, and recording from various origins, such as observational
194 studies, controlled experiments, publications, and databases. The collected data should be pertinent
195 to the problem, accurate, and suitable for ML model development. Along with data collection, one
196 needs to document the sources, methods, and potential biases associated with the data to ensure
197 transparency and reproducibility.

198

199 **3.2.3 Data preprocessing**

200 Data preprocessing is critical in a ML workflow that involves cleaning (i.e., filtering),
201 transforming, and organizing raw data. Cleaning is identifying and rectifying errors, inaccuracies,
202 inconsistencies, and anomalies in the collected data, which usually involves filtering missing
203 values, detecting outliers, removing duplicate entries, and converting data types. Normalization,
204 as an example of data transformation, is scaling data to fit within a specific range, such as 0 to 1.
205 Normalization equalizes the contribution of each feature to the model training and avoids any
206 single attribute disproportionately influencing the results. Organizing raw data entails structuring
207 and arranging the data in a manner that is optimal for the specific ML algorithms to be applied.

208

209 **3.2.4 Model training**

210 The first step in data training is to divide the dataset into training, validation, and test datasets
211 (Hastie et al., 2009). The training dataset trains the model, the validation dataset tunes
212 hyperparameters, and the test dataset evaluates the performance of the model. The second step is
213 to select an appropriate model on the basis of the nature of the problem. Not all algorithms work
214 equally well for all types of problems. During training, the model recognizes the relationship
215 between the inputs and outputs and minimizes the difference between predicted and actual outputs
216 by iteratively adjusting the model parameters.

217

218 **3.2.5 Model evaluation**

219 Evaluation metrics differ on the basis of the nature of the problem (Xie et al., 2023). Evaluation
220 metrics assess the performance of classification models, revealing the effectiveness of the
221 classification models and their ability to distinguish between classes. Accuracy refers to the
222 proportion of correctly classified instances to the total instances. Precision measures the proportion
223 of correctly predicted positive instances among all predicted positives. Recall gauges the
224 proportion of correctly predicted actual positive instances. The F1-score combines precision and
225 recall into a single metric, offering a balanced view of the accuracy of a model (Sokolova and
226 Lapalme, 2009). The confusion matrix provides a tabular representation of true positive (TP), true
227 negative (TN), false positive (FP), and false negative (FN) predictions (James et al., 2013).
228 Receiver operating characteristic (ROC) curves illustrate the trade-off between the TP rate and the
229 FP rate at different classification thresholds with the area under the curve (AUC) summarizing the
230 performance of the curve (Bradley, 1997). In regression tasks, the mean squared error (MSE) and

231 root mean squared error (RMSE) quantify the average squared differences between predicted and
232 actual values, and the mean absolute error (MAE) measures the average absolute differences
233 (Willmott and Matsuura, 2005). Additionally, the coefficient of determination (R^2) indicates the
234 proportion of variance in the target variable explained by the model (Steel and Torrie, 1960).

235

236 **4. Machine learning to ensure safe drinking water supply: The physical perspective**

237 **4.1. To predict and manage drinking water production and demand**

238 ML is promising in ensuring safe drinking water supply from the physical perspective (Table 1).
239 Insufficient production of drinking water could constrain socioeconomic development and
240 population growth (Grey and Sadoff, 2007). A hybrid model that combines genetic algorithms (GA)
241 and GA artificial neural networks (GA-ANN) can predict drinking water production (Figure 2a)
242 (Zhang et al., 2019). The model uses predictors such as temperature, chemical oxygen demand
243 (COD), and electricity and chemical consumption. GA can optimize weights and biases, enhancing
244 the prediction accuracy of the relationship between the inputs and outputs of the ML model. GA
245 can also increase the tolerance to imprecisions, uncertainties, and approximations in the inputs.
246 The hybrid GA-ANN was trained and validated with monthly data from 45 water utilities across
247 China. The performance of GA-ANN ($R^2 = 0.93$) is substantially better than a three-layered ANN
248 ($R^2 = 0.71$) when more training data are incorporated. The GA-ANN effectively forecasts
249 fluctuations in water production for various scenarios, highlighting its feasibility in promptly and
250 appropriately adjusting water treatment operations. To assess drinking water demand patterns, one
251 can apply unsupervised learning algorithms to raw time-series data of drinking water consumption.
252 For instance, the hierarchical K-means algorithm can classify drinking water consumption patterns
253 (Leitão et al., 2019). In that algorithm, daily time-series demand with hourly records served as the

254 sample inputs within a 24-dimensional feature space to identify dense and distinctly separated
255 temporal patterns on drinking water demand. In contrast to directly clustering the drinking water
256 demand patterns, short-time water demand forecasting is more intriguing and has more practical
257 merits in optimizing DWDSs. A gated recurrent unit network (GRUN) predicts short-term water
258 demand for the upcoming 15 min and 24 h using a time step of 15 min (Guo et al., 2018). The
259 GRUN model has a higher accuracy with a mean absolute percentage error (MAPE) of 2.06% than
260 the conventional three-dense-layered ANN (MAPE = 2.46%) and the seasonal autoregressive
261 integrated moving average (SARIMA) model (MAPE = 2.57%) for the 15-min prediction. For the
262 24-h prediction, the GRUN model also achieves more precise forecasts with MAPE ranging from
263 4.33% to 4.96%. Another study used three ML models to forecast both short-term and long-term
264 water demand encompassing daily, weekly, and monthly intervals in Tehran, Iran (Ghiassi et al.,
265 2017). The models include a dynamic artificial neural network (DAN2), a focused time-delay
266 neural network (FTDNN), and a KNN. DAN2 based on the adaptive network principle allows the
267 architecture to adjust dynamically responding to data-driven learning. Given its inherent design,
268 DAN2 is promising at time series forecasting, catering to datasets characterized by evolving
269 temporal patterns (Ghiassi et al., 2005). For instance, DAN2 achieves remarkable prediction
270 accuracies (96% for daily, 99% for weekly, and 98% for monthly water demand forecasts) and
271 outperforms both FTDNN and KNN (Ghiassi et al., 2017).

272

273 **4.2 To monitor and predict drinking water pipeline failures**

274 Pipeline failures in engineered water systems cause significant water loss and contaminate
275 municipal water (Renwick et al., 2019). These failures can introduce microbes and chemicals from
276 the surroundings into distributed water. A novel burst location identification framework by fully-

277 linear DenseNet (BLIFF) framework can detect pipe burst locals (Figure 2b) (Zhou et al., 2019).
278 BLIFF relies on deep learning through the fully-linear DenseNet (FL-DenseNet) model. BLIFF
279 supplants the convolutional layers in DenseNet with linear connections and omits pooling layers.
280 With the real-time pressure measurements as the inputs, BLIFF generates the likelihood values of
281 a burst for each pipe in the potential burst district. The prediction accuracies (62.35% to 98.58%)
282 of BLIFF are almost two times the original DenseNet model. This remarkable improvement is due
283 to the enhanced ability of the linear-connection layer to discern global features within the pressure
284 signals. The effective deployment of deep learning methods such as BLIFF corroborates the
285 viability of pressure values in burst localization, countering prior assertions of their insensitivity
286 to burst events (Bakker et al., 2014; Mounce et al., 2010). Another cutting-edge paradigm,
287 advanced meta-learning (AdvaML), can predict the failure of drinking water pipelines (Almheiri
288 et al., 2021). AdvaML is also a deep learning model but is rooted in the ANN architecture. AdvaML
289 comprises an input layer with 33 neurons (mirroring the 33 predictors including pipe and climate
290 data), four hidden layers, and an output layer that yields the failure/hazard index of a pipe. AdvaML
291 not only forecasts the risk index associated with pipe failure but also detects pivotal determinants
292 of pipeline failure timing. Of these determinants, the number of lanes and chlorine residual
293 concentration are paramount, collectively contributing approximately 9% to the lifespan analysis
294 of DWDSs. Notably, AdvaML tackles the challenge of limited data, promoting the applications of
295 deep learning to assess pipe failure. Compared with alternative models like cox-proportional
296 hazards (Cox-PH), survival support vector machine (SSVM), and random survival forest (SRF),
297 AdvaML has commendable performance even with scant training data because of its knowledge
298 transfer from initial parameterization to the ultimate learning phase.
299

300 While inherent system vulnerabilities cause pipeline failures, external factors exacerbate the issue
301 (Fan et al., 2023; Folkman, 2018). Climatic extremes and weather disasters, such as wildfires,
302 become more frequent because of climate change and threaten drinking water infrastructures. In
303 response, researchers leverage ML to better understand and predict the impact of these disasters
304 on water pipes. Two ensemble ML models (RF and XGB) can assess and forecast the repercussions
305 of calamities on water treatment infrastructures (Park et al., 2020). These models incorporate 23
306 variables encompassing facility specifications and operational data from 419 water utilities in
307 South Korea. The models project the total disaster index (TDI), a metric signifying the effects or
308 damages wrought by three predominant disasters (typhoons, heavy rainfalls, and earthquakes) on
309 DWDSs. While both RF and XGB have commendable predictive prowess concerning the TDI,
310 XGB slightly outperforms in most scenarios. Another study developed four models, a linear
311 regression-based repair rate (RR) method, LR, boosted regression trees (BRT), and RF, to predict
312 pipeline damage during an earthquake (Bagriacik et al., 2018). The models incorporate parameters
313 such as ground shaking, permanent ground deformation, pipe material, pipe diameter, year of
314 installation, type, and trench backfill type. Each model demonstrates unique strengths. The BRT
315 model has the best overall predictive performance, while the LR model is instrumental in
316 highlighting the influence of pipe materials and trench types on pipeline damage.

317

318 **5. Machine learning to monitor and ensure chemical drinking water quality**

319 **5.1. To assess and control disinfection by-products in municipal water**

320 **5.1.1. To predict the concentration of disinfection by-products in drinking water**

321 ML is useful in monitoring and ensuring chemical drinking water quality (Table 2). Drinking water
322 disinfection is critical to ensuring microbial drinking water quality and safeguarding public health

323 (Zhang and Lu, 2021c). Disinfection is effective in killing pathogens, impeding microbiological
324 recontamination, and inhibiting biofilm development in drinking water (Mazhar et al., 2020).
325 Chlorine-based disinfectants, such as free chlorine (e.g., chlorine gas and sodium hypochlorite),
326 bound or combined chlorine (e.g., monochloramine), and chlorine dioxide, are widely used in
327 water treatment because of their cost-effectiveness and high efficiency (Gupta and Ali, 2013) (Jefri
328 et al., 2022; Zhang et al., 2018). Nonetheless, when these disinfectants interact with natural organic
329 matter (NOM) and anthropogenic compounds (such as pharmaceuticals and antimicrobials), they
330 generate DBPs such as trihalomethanes (THMs), haloacetic acids (HAAs), haloketones (HKs),
331 haloacetonitriles, halophenols, and halopropanoles (Li et al., 2023a) (Favere et al., 2021; Xiao et
332 al., 2023). DBPs are harmful and even carcinogenic with significant health risks (Pandian et al.,
333 2022; Zhou et al., 2023). Therefore, monitoring and controlling DBPs in drinking water is vital to
334 public health (He et al., 2021; Helte et al., 2023; Redondo-Hasselerharm et al., 2022). Conventional
335 methods for monitoring DBPs require expensive instrumentation such as gas chromatography (GC)
336 and liquid chromatography (LC) combined with mass spectrometry (MS) and complicated pre-
337 treatment processes. Thus, those conventional methods are labor-intensive, costly, and time-
338 consuming, limiting the ability of water utilities to reduce DBP formation. By contrast, ML models
339 for monitoring DBPs in drinking water are accurate, efficient, inexpensive, and easy to handle
340 (Balogun et al., 2021; Jia et al., 2021; Podgorski and Berg, 2022).

341
342 Finding the optimal disinfectant dosages to simultaneously minimize the formation of DBPs and
343 pathogen regrowth in finished water is crucial (He et al., 2021; Zhang and Lu, 2021c). Nevertheless,
344 reaching this goal with traditional methods is time-consuming, expensive, and complex. Therefore,
345 ML stands out as it is effective in predicting the formation of DBPs, significantly reducing capital

346 and human investment. Common input parameters in these models are operational parameters
347 (temperature and time) and water quality variables, such as pH, UV_{254} , and the concentrations of
348 dissolved organic carbon (DOC), chloride (C_{resCl^-}), bromide (C_{Br^-}), nitrite nitrogen ($C_{NO_2^-N}$), and
349 ammonium nitrogen ($C_{NH_4^+-N}$) (Deng et al., 2021; Hong et al., 2020; Hu et al., 2023; Lin et al.,
350 2020; Pan et al., 2023; Singh and Gupta, 2012). The outputs are the concentrations of DBPs, such
351 as THMs, HAAs, and HKs
352
353 A study developed three ML models [ANN, SVM, and gene expression programming (GEP)] to
354 forecast the formation of DBPs on the basis of a 63-point dataset (Singh and Gupta, 2012).
355 Specifically, pH, temperature, contact time (t), the concentration of bromide, and DOC-normalized
356 chlorine dose (Cl_2/DOC) are the inputs. The concentration of THMs in chlorinated river water is
357 the output. SVM outperforms the other two models, exhibiting the highest R^2 and the lowest RMSE
358 values. Furthermore, sensitivity analysis reveals that pH, HRT, and temperature are the top three
359 contributors to DBP formation. In addition, radial basis function (RBF) based ANN models can
360 predict the formation of common DBPs such as HAAs (Lin et al., 2020), THMs (Hong et al., 2020),
361 and HKs (Deng et al., 2021) in DWDSs. For instance, a study extracted 64 representative data
362 points from the literature to predict HAA formation using pH, temperature, DOC, UV_{254} , C_{resCl^-} ,
363 C_{Br^-} , $C_{NO_2^-N}$, and $C_{NH_4^+-N}$ as the inputs (Figure 3a) (Lin et al., 2020). The RBF ANN model
364 outperforms the linear and log-linear models by 21% to 47 % in accuracy, respectively. Therefore,
365 the RBF ANN model is promising in assessing DBP formation and optimizing disinfection. A
366 follow-up study used 64 data points to train an RBF ANN model to predict THM formation (Hong
367 et al., 2020). The RBF ANN model achieves accuracies between 92% and 98% and regression

368 coefficients ranging from 0.76 to 0.93, outperforming the linear and log-linear models and
369 demonstrating its superiority to uncover complex non-linear patterns in THM formation. Even
370 when trained with fewer water quality variables, a fusion of grey relation analysis with RBF ANN
371 could provide superior prediction results. Furthermore, an RBF ANN model trained with 63 data
372 points from tap water predicts the formation of HK (Deng et al., 2021). Both RBF ANN and back
373 propagation (BP) ANN models outstrip the linear and log-linear models with RBF ANN displaying
374 higher accuracies in both internal and external validations. Another study explored the application
375 of a decision tree boost (DTB) model to predict the concentrations of THM4 and HAAs (Pan et
376 al., 2023). The study examined correlations between water quality parameters and mixed
377 chlorine/chloramine species. The work then selected seven variables such as NH_2Cl , NHCl_2 plus
378 organic chloramines, pH, total dissolved nitrogen (TDN), nitrite, total organic carbon (TOC), and
379 NH_4^+ as the independent variables to predict THM4 and HAAs. The DTB model demonstrates
380 enhanced prediction accuracy with R^2 values of 0.56 for THM4 and 0.65 for HAAs, while the
381 inclusion of organic chloramines improves the prediction precision. Additionally, a study
382 implemented multiple ML models to predict emerging DBPs in small DWDSs across Canada by
383 analyzing data from eleven such networks (Hu et al., 2023). The models use parameters like
384 temperature, total chlorine residual, DOC, turbidity, pH, conductivity, and UV_{254} to predict the
385 concentrations of THMs, HAAs, dichloroacetonitrile (DCAN), chloropicrin (CPK), and
386 trichloropropanone (TCP). Among the evaluated models, support vector regression (SVR) and
387 Gaussian process regression (GPR) show superior performance with SVR exhibiting the highest
388 prediction accuracy ($R^2 = 0.94$) and stability for DCAN and TCP, while GPR is optimal for
389 predicting CPK ($R^2 = 0.92$).

390

391 UV-visible and fluorescence spectral techniques are preferred to monitor DBPs (Krasner et al.,
392 2006; Rodriguez et al., 2004). Fluorescence spectroscopy is sensitive to assessing the
393 characteristics and reactivity of NOM owing to its minimal sample preparation requirement and
394 short acquisition time (Pifer and Fairey, 2012). However, the complex high-dimensional
395 characteristics of fluorescence spectroscopy make it difficult to predict DBP formation. The
396 resources and time required for DBP analysis restrict the capacity of water utilities to reduce DBP
397 formation. This situation has prompted the development of ML models that can predict DBP
398 formation. Since DBPs form from the reaction between chlorine-based disinfectants and NOM,
399 including NOM measurement in these models is critical. Nonetheless, NOM is a diverse group of
400 organic molecules with complex characteristics, adding to the complexity of capturing the
401 reactivity of NOM with chlorine-based disinfectants (Wagner and Plewa, 2017).

402
403 Autoencoder-neural networks (AE-NN) can predict the concentrations of both THMs and HAAs
404 in river water from fluorescence spectra (Figure 3b) (Peleato et al., 2018). To manage the high
405 dimensionality of fluorescence spectra, the researcher applied three dimension-reduction
406 techniques [AE-NN, parallel factors analysis (PARAFAC), and PCA] before further analysis.
407 Afterward, the researcher trained NN to identify fluorescence regions associated with DBP
408 formation and to predict DBP concentrations. The AE-NN model has superior predictive
409 accuracies for THMs and HAAs, achieving validation MAE values of 9.65 $\mu\text{g/L}$ and 9.64 $\mu\text{g/L}$,
410 respectively. These figures exceed the performance of PCA, which has higher errors of 13.19 $\mu\text{g/L}$
411 for THMs and 11.92 $\mu\text{g/L}$ for HAAs. Furthermore, the precision of the AE-NN model surpasses
412 that of parallel factors analysis, which has MEA values of 20.39 $\mu\text{g/L}$ for THMs and 14.00 $\mu\text{g/L}$
413 for HAAs. In addition, convolutional neural networks (CNNs) can predict DBP concentrations

414 from fluorescence spectra without extensive data pre-processing (Peleato, 2022). Compared with
415 multilayer perceptron (MLP) and dimensionality reduction techniques, CNNs not only exhibit
416 superior prediction accuracies for THMs and HAAs but also identify the fluorescence spectra
417 regions highly associated with DBP formation.

418

419 **5.1.2. To elucidate the formation mechanisms of disinfectant by-products in drinking water**

420 ML is promising in predicting the formation of DBPs using either water quality metrics and
421 operational parameters or via online spectrum monitoring. However, even with the knowledge of
422 DBP concentrations, their removal from drinking water remains costly and inefficient (Bond et al.,
423 2011; Rodriguez et al., 2004). An effective approach for DBP control is to remove DBP precursors
424 and prevent them from reaching the clear wells in water utilities (Bond et al., 2012; Krasner et al.,
425 2013). This needs a comprehensive understanding of the mechanisms for the formation of DBPs.

426

427 A multiple linear regression (MLR) model can predict the production of chloroform (a THM
428 compound) from organic precursors (Figure 3c) (Bond and Graham, 2017). Relying on 211
429 precursors from 22 studies, the MLR model uses 19 descriptors as the inputs and chloroform yield
430 as the output. The well-trained MLR model has a promising prediction accuracy with an R^2 value
431 of 0.91 and an RMSE value of 8.93 mol/mol. Further chemical insights pinpoint that functional
432 groups, such as hydroxyl, chlorine, and carboxyl groups, significantly affect chloroform formation.
433 ML can also forecast the formation of HAAs from the interaction between organic precursors and
434 free chlorine (Cordero et al., 2021). The training dataset comprises 283 organic compounds and
435 732 chemical descriptors as the inputs with HAA yield as the output. These organic compounds
436 are converted into 2D and 3D chemical descriptors with their SMILE strings used for ML

437 compatibility. Three common ML models (RF, SVR, and MLP) are selected because they can
438 handle nonlinear problems, activity cliffs, and high dimensions in addition to MLR as a benchmark.
439 The RF model is the top performer with the lowest RMSE values of 1.05 and 1.19 for
440 dichloroacetic acid (DCAA) and trichloroacetic acid (TCAA), respectively. The crucial predictors
441 of TCAA formation are the number of aromatic bonds, hydrophilicity, and electrotopological
442 descriptors related to electrostatic interactions and the atomic distribution of electronegativity.

443

444 **5.1.3. To investigate alternative disinfectants to reduce disinfectant by-product formation**

445 Since chlorine-based disinfectants produce harmful DBPs, alternative disinfectants for drinking
446 water disinfection attract attention. An alternative disinfectant is ozone. Unlike chlorination,
447 ozonation does not produce chlorinated THMs and HAAs. Ozone, therefore, provides a two-fold
448 benefit: it is effective and does not generate chlorinated DBPs. However, ozonation produces
449 various other DBPs (Mao et al., 2014; Richardson et al., 1999). The occurrence and toxicity of
450 these ozonated DBPs are a concern (Simpson and Mitch, 2022; Srivastav et al., 2020). Therefore,
451 while ozone does not generate chlorinated DBPs, its use necessitates careful consideration of the
452 potential for other toxic byproducts. In this section, we review how ML models have been
453 developed and applied in the field of ozonation. ML can help control the formation of bromate and
454 reduce micropollutants, microbial indicators, and organic contaminants during ozonation.

455

456 ANN has two advantages when compared with MLR in controlling bromate formation during
457 ozonation (Legube et al., 2004): First, ANN with an R^2 value of 0.98 is more accurate than MLR.
458 Second, ANN classifies model variables (predictors) in the descending order of impact: ozone dose,
459 $C_{\text{NH}_4^+-\text{N}}$, C_{Br^-} , pH, temperature, DOC, and alkalinity. While ANN has superior performance, the

460 simplicity of MLR is attractive. However, a key limitation of MLR is that its accuracy decreases
461 with an increased number of samples because MLR cannot effectively process nonlinear
462 components.

463
464 ML models based on routinely measured physical-chemical water quality parameters can also
465 predict the changes in micropollutant concentration during ozonation. For instance, RF models can
466 predict the reduction of micropollutants during ozonation(Cha et al., 2021) (Figure 3d). That study
467 introduced four distinct RF models, each incorporating standard predictors such as pH, alkalinity,
468 and DOC. These models have unique inclusions of fluorescence excitation-emission matrix
469 (FEEM) data at different resolutions. These models are as FEEM-Free, FEEM-LowRes, FEEM-
470 HighRes, and FEEM-FullRes, each varying in the resolution of FEEM data used as the unique
471 predictors. A comparison of these four models shows that integrating FEEM data results in more
472 accurate predictions of ozone exposures. The high-resolution FEEM data yield better predictions
473 for microplastic abatement ($R^2 = 0.904$; RMSE = 6.6%). However, the improvement in prediction
474 accuracy when using FEEM data is less substantial for predicting microplastic abatement than for
475 predicting the oxidant exposures.

476
477 Machine-learning-quantitative-structure-property-relationship (ML-QSPR) methods can calculate
478 the rate constant (k_{O_3}) of the reactions between ozone and micropollutants (Gupta and Basant,
479 2016; Huang et al., 2020; Shi et al., 2022; Sudhakaran and Amy, 2013). Generally, nonlinear
480 models outperform their linear counterparts. For instance, an MLR method (Sudhakaran and Amy,
481 2013) and an SVM method (Huang et al., 2020) have performance levels with R^2 values of greater
482 than 0.75 and 0.78, respectively. Conversely, a DTB model has a notably higher R^2 value of greater

483 than 0.97 (Gupta and Basant, 2016). A recent study thoroughly compared several models including
484 MLR, SVM, DT, RF, and deep neural network (DNN) for predicting $\log k_{O_3}$ (Shi et al., 2022). Of
485 these, the RF model has the highest effectiveness with a peak R^2 value of 0.91. The RF model has
486 two primary benefits: robustness and a lower tendency toward overfitting. On the other hand, the
487 DT model is a complex increase in its structure, subsequently raising the overfitting risk. In
488 addition, DNN, promising at recognizing nonlinear features, underperforms in predicting $\log k_{O_3}$.
489 A similar situation occurs when ML models predict the elimination of recalcitrant trace organic
490 compounds by ozonation for municipal wastewater reuse (Park et al., 2015). Specifically, ANN
491 was susceptible to overfitting issues. Incorporating PCA into ANN creates a PC-ANN workflow,
492 which addressed these issues. PCA transforms the input variables to linearly independent variables,
493 thereby resolving the issue of collinearity among explanatory variables. The PC-ANN method (R^2
494 = 0.934) surpasses the standalone ANN ($R^2 = 0.914$) in terms of predictive power.

495

496 **5.2. To monitor and control heavy metals in drinking water**

497 **5.2.1. To monitor heavy metals in drinking water**

498 The quality of drinking water in engineered water systems can deteriorate because of the
499 detachment of heavy metals from the water pipes into the bulk water (i.e., leaching) (Mays, 2000;
500 Proctor et al., 2020). The detrimental effects of heavy metal toxicity on all living beings have been
501 well-documented (Umesh C. Gupta, 2011; Valko et al., 2005). Therefore, this section focuses on
502 the applications of ML in assessing heavy metals in drinking water.

503

504 Several studies used ML to estimate the concentration of heavy metals in drinking water. For
505 instance, a continuous on-site, *in situ* system can estimate lead (Pb) concentrations in municipal

506 water (Oh et al., 2021). The system leverages the SVR algorithm, supplanting traditional
507 mathematical models confined to analyzing stationary ions in a solid substrate. By using the radio-
508 frequency reflection coefficient of the raw trace data, the system predicts Pb ion concentration
509 with a resolution of 1 $\mu\text{g/L}$ and an RMS prediction error of 0.71 $\mu\text{g/L}$ in the presence of interfering
510 metal ions such as Cu^{2+} , Fe^{3+} , and Zn^{2+} . Other than estimating heavy metal concentrations in
511 individual samples, ML is promising in broader analytical applications. For instance, ML is useful
512 in the spatial interpolation of environmental variables, significantly enhancing its performance (Li
513 et al., 2011). This approach has been substantiated in developing spatial interpolation maps
514 depicting the concentrations of heavy metals such as Fe, Mn, Ni, Pb, and Zn in groundwater (i.e.,
515 a water source of drinking water) (De Jesus et al., 2021). That study combined ML and
516 geostatistical interpolation (MLGI) to leverage an ANN-based algorithm to augment the efficacy
517 and robustness of the spatial interpolation mapping. Furthermore, the MLGI approach
518 comprehensively assesses the carcinogenic risks of heavy metals through *in situ* measurements.
519 The approach produces detailed spatial maps delineating metal concentrations and estimates health
520 quotient indices (HQI) to offer a more refined risk assessment (Senoro et al., 2022). While the
521 integration of ML algorithms elevates the efficacy and robustness of spatial interpolation,
522 traditional interpolation techniques still have a significant role in this domain. For instance, a
523 spherical semi-variogram model relying on the classic Kriging interpolation technique can monitor
524 the temporospatial distribution of residual aluminum (Al) in a DWDS (Figure 4a), thereby
525 highlighting the enduring relevance and applicability of traditional interpolation methods (Tian et
526 al., 2020).

527

528 5.2.2. To control heavy metals in drinking water

529 Adsorption is proficient in mitigating heavy metal contamination in drinking water (Joseph et al.,
530 2019; Wołowiec et al., 2019). The removal of heavy metals by adsorption has highly stochastic,
531 non-linear, and non-stationary dynamics coupled with redundancy (Bhagat et al., 2020). Many ML
532 techniques can enhance the precision and efficacy of predicting heavy metal adsorption dynamics.
533 Common predictors for those ML techniques are adsorbent dosage, operating temperature, contact
534 time, and pH, whereas the output is the removal efficiency of heavy metals. Other variables can
535 also be incorporated into the predictive models such as the initial concentration of heavy metals,
536 the specific surface area of metal-organic frameworks (MOFs), and the presence of anions (Abdi
537 and Mazloom, 2022). That study predicted the adsorptive removal of arsenate [As(V)] using four
538 ML techniques: light gradient-boosting machine (LightGBM), XGB, gradient-boosted decision
539 trees (GBDT), and RF. In addition, each predictive model can predict the removal efficiencies of
540 multiple heavy metals. For instance, ANN can predict the removal efficiencies of Al, Cd, Co, Cu,
541 Fe, and Pb (Hamidian et al., 2019). The LightGBM model has the most precise predictions for
542 As(V) adsorption (Abdi and Mazloom, 2022). However, deep-layer-structured ANN has a higher
543 accuracy particularly when employing radial basis functions as activation functions in the hidden
544 layer (Hamidian et al., 2019).

545

546 **5.3. To monitor nitrification in engineered water systems**

547 Nitrification is a serious issue in chloraminated engineered water systems (Sathasivan et al., 2008;
548 Shi et al., 2020). During nitrification, ammonia-oxidizing microbes oxidize free ammonia to nitrite,
549 and nitrite-oxidizing bacteria further oxidize nitrite to nitrate. Free ammonia thus initiates
550 nitrification. Free ammonia appears in chloraminated engineered water systems because of
551 excessive free ammonia dosing at water utilities, chloramine decay within water pipes, and the

552 reactions between chloramines and reducing agents. Nitrification deteriorates water quality by
553 destroying chloramine residuals, releasing free ammonia, promoting microbial (re)growth, and
554 producing nitrite and nitrate. Therefore, monitoring nitrification in engineered water systems is
555 critical to ensuring drinking water quality and protecting public health.

556
557 ML is useful in investigating drinking water nitrification. A supervised ML technique, the NB
558 classifier, relies on biomass and microbiome datasets to detect nitrification (Gomez-Alvarez and
559 Revetta, 2020). After being trained with microbiome indicators, the model has a binary
560 classification accuracy of up to 85% with an AUC of 0.825 when distinguishing between
561 nitrification and stable events. An SVR (Figure 4b) is trained with nitrate/nitrite spectra at various
562 wavelengths (predictor variables) along with corresponding nitrate/nitrite concentrations (response
563 variables) (Hossain et al., 2021). The SVR model negates the need for any chemical supplements,
564 is easy to use, and can reach a high level of precision of up to ± 0.01 mg N/L.

565 566 **6. Machine learning to monitor and ensure microbiological drinking water quality**

567 **6.1. To monitor and mitigate opportunistic pathogens in municipal water with a focus on** 568 ***Legionella***

569 ML is useful in monitoring and ensuring microbial water quality in engineered water systems
570 (Table 3). The microbial community within engineered water systems is diverse, encompassing
571 general heterotrophic bacteria, protozoans (such as amoebae, ciliates, and slime molds), and OPs.
572 Dominant water-related OPs are *Legionella* (especially *L. pneumophila*), *Mycobacterium* (e.g.,
573 nontuberculosis mycobacteria or NTM and *M. avium* complex or MAC), *Pseudomonas aeruginosa*,
574 *Vermamoeba vermiformis*, *Naegleria fowleri*, and *Acanthamoeba* (Donohue et al., 2019; Isaac and

575 Sherchan, 2020; Lytle et al., 2021). In municipal water, OPs are the most significant parameter for
576 microbial drinking water quality because of their frequent occurrence, high concentrations, high
577 resistance to disinfectant residuals, and association with drinking-water-related diseases such as
578 Legionnaires' disease (Zhang and Lu, 2021a). Therefore, closely monitoring OPs in engineered
579 water systems is critical to assessing drinking water quality and protecting public health.
580 *Legionella* is the most important OP in municipal water. In addition, compared with conventional
581 microbial drinking water quality indicators such as fecal coliforms, *Legionella* is a better candidate
582 for indicating microbial drinking water quality (Zhang and Lu, 2021a). Therefore, in this section,
583 we focus on the applications of ML in monitoring and controlling *Legionella* in municipal water.
584

585 Studies using ML to assess the risks of OPs remain limited. An early work mitigated the
586 proliferation of *Legionella* in premise plumbing by controlling environmental variables (Sincak et
587 al., 2014). Using water flow and temperature as the inputs, that study presented a NN-based
588 simulator designed using an approximate reasoning architecture (NARA) neuro-fuzzy system to
589 predict and simulate water tank temperature profiles. The simulator emulates conditions that
590 inhibit the spread of *Legionella* in water networks. The NARA-based simulator achieves a high
591 fidelity in mimicking water tank temperatures with an accuracy exceeding 97%. Another study
592 integrated both unsupervised and supervised ML techniques to correlate the spread of *Legionella*
593 with environmental variables in retirement homes, health-related facilities, tourism-related
594 buildings, and swimming-pool environments in Italy (Brunello et al., 2022). That study used an
595 unsupervised learning algorithm to identify the spatiotemporal distribution of atypical *Legionella*
596 through an ordinal regression model. The results indicate how the distribution is correlated with
597 the healthcare facilities. Hospitals had the highest contamination cluster locations, indicating a

598 strong correlation between the propagation of *Legionella* and both the nature of the facilities and
599 broader geographical characteristics. That work also used supervised ML to assess the serotypes
600 of *Legionella* and to anticipate the corresponding contamination levels. For serogroup assessment,
601 XGBoost, LR, and SVM Classifier were used and compared. XGBoost shows superior
602 performance with an overall classification accuracy of 0.71. The contribution of each predictor to
603 the final classification was evaluated by the Shapley values, which quantify the contribution of
604 each variable to the outputs of a ML model by comparing the effect of the outputs relative to the
605 average across all inputs. The geographical location of a sample is the most important parameter
606 but is useful only when combined with other predictors. For contamination level prediction, all
607 three models demonstrate low performance with the highest accuracy of 0.57 from XGBoost.

608

609 **6.2. To detect *Cryptosporidium* and *Giardia* in drinking water**

610 *Cryptosporidium* and *Giardia* are protozoan parasites in municipal water with substantial public
611 health risks by causing cryptosporidiosis and giardiasis, respectively (CDC, 2021a; b). These
612 pathogens are notably resilient to disinfectants such as chlorine, challenging water treatment
613 (Adeyemo et al., 2019). Therefore, detecting and controlling *Cryptosporidium* and *Giardia* is
614 critical to maintaining drinking water quality. In this section, we discuss the performance of these
615 ML models in detecting *Cryptosporidium* and *Giardia*.

616

617 ML to detect and analyze *Cryptosporidium* and *Giardia* has high robustness and precision. For
618 instance, deep-learning-based image classification models such as ParasNet (Xu et al., 2020) and
619 MCellNet (Luo et al., 2021) are accurate in detecting these two parasites in drinking water. They
620 show the power of ML in classifying parasites from the cell-level scattering images. In addition, a

621 linear ML model can predict the contamination of these two parasites (Figure 5a) (Ligda et al.,
622 2020), offering a valuable tool to control waterborne diseases.

623
624 ParasNet uses an eight-layer CNN to determine whether particles in cell-level scattering images
625 from drinking water are *Cryptosporidium* and *Giardia*. The model has superior performance
626 compared with a traditional handcraft SVM regarding both detection accuracy and processing
627 speed. For instance, ParasNet can reach above 95.6% detection accuracy with analysis speeds of
628 up to 100 frame-per-second (fps). MCellNet, another image classification pipeline, uses a DNN
629 optimized from MobileNetV2 to recognize objects (Sandler et al., 2018). MCellNet includes a
630 convolutional layer, six inverted residual blocks (IRBs), a flattened layer, and a fully connected
631 layer. MCellNet can process images from flow cytometry to classify *Cryptosporidium* and *Giardia*.
632 Compared with ParasNet, MCellNet achieves a higher detection accuracy of above 99.6% with a
633 346-fps analysis speed. The superior accuracy and expedited analysis of MCellNet are due to the
634 cascading six IRBs in the model.

635
636 An alternative statistical model uses linear discriminant function analysis (LDFA) to predict the
637 contamination of *Cryptosporidium* and *Giardia* in drinking water (Ligda et al., 2020). That model
638 uses microbiological, physicochemical, and meteorological parameters as predictors to classify the
639 contamination of *Cryptosporidium* and *Giardia* into four categories: none, low, moderate, and high
640 (oo)cysts concentrations. LDFA has accuracies of 75% and 69% in predicting the contamination
641 of *Cryptosporidium* and *Giardia*, respectively.

642

643 **6.3. To assess biofilm development in engineered water systems**

644 Research on the applications of ML to study biofilm development in engineered water systems is
645 scarce. ML models in this research area use physical (such as hydraulic) factors to predict the
646 dynamics of biofilm development, where the heterotrophic plate count (HPC) is the common
647 output. Established ML algorithms are preferred models to study the dynamics of biofilm
648 development such as NB, RT, and RF (Ramos-Martínez et al., 2014; 2016). These algorithms have
649 high prediction accuracy and provide a deeper understanding of the impact of physical factors on
650 biofilm development in engineered water systems. For instance, a Bagging naïve Bayesian tree
651 (B-NBT) model proposes optimal flow velocities for different types of pipes to mitigate biofilm
652 development (Ramos-Martínez et al., 2014). To control biofilm accumulation in DWDSs, water
653 utilities may avoid cement pipes, implement medium- or high-flow velocities in metal pipes and
654 sustain water ages above 0.035 in plastic pipes. The 'water age' is a synthetic index derived from
655 the normalized hydraulic retention time (HRT) and the distance from the disinfection source.

656

657 Recent studies have enabled more detailed, single-cell level analyses and predictions of biofilm
658 development in engineered water systems (Berne et al., 2018). In addition, researchers have
659 expanded algorithms to incorporate deep learning methods (Jelli et al., 2023; Weigert et al., 2020).
660 These innovative approaches can enhance the understanding of biofilm dynamics in various
661 settings including engineered water systems. A study developed a refined deep learning technique
662 for the single-cell segmentation of 3D bacterial biofilms and made two breakthroughs (Figure 5b)
663 (Jelli et al., 2023). First, it expedited the annotation process of bacterial cells within 3D imagery,
664 enhancing the efficiency of data analysis. Second, it optimized the application of StarDist (Weigert
665 et al., 2020), a cutting-edge CNN-based algorithm. The refined deep learning technique achieves
666 unprecedented accuracy in biofilm segmentation, surpassing other algorithms under scrutiny.

667 Moreover, the technique tracks cell lineages and enables precise measurements of bacterial growth,
668 offering unprecedented insights into biofilm structures and development. While the application of
669 ML to assess biofilm development in engineered water systems remains to be explored, such
670 advancements herald a promising avenue for future research of biofilm development in drinking
671 water.

672

673 **6.4. To predict antimicrobial resistance risks and track the sources of antimicrobial** 674 **resistance genes in drinking water**

675 **6.4.1. To predict antimicrobial resistance risks in drinking water**

676 Antimicrobials have been extensively used since the 1920s in the medical industry, animal
677 husbandry, and other fields (Chang et al., 2015; Hutchings et al., 2019; Prescott, 2017; WHO,
678 2021). More than 80% of antimicrobials used in humans and animals are not metabolized and
679 excreted from the bodies. Antimicrobials cause the issue of antimicrobial resistant bacteria (ARB)
680 and ARGs. ARB and ARGs can enter drinking water and cause antimicrobial resistance (AMR),
681 significantly threatening public health (Roca et al., 2015; Walesch et al., 2023).

682

683 Detecting ARB and ARGs with conventional assays is challenging. Therefore, estimating the risks
684 of AMR in drinking water is time-consuming and costly. A novel ML approach predicting the risks
685 of AMR can overcome this issue (Wu et al., 2022). That approach maps the relative risk scores of
686 AMR from continuous values into binary (0 and 1) labels using a predefined threshold (Goh et al.,
687 2022). The study formulated the relative risks of AMR as a binary classification task and
688 investigated multiple ML models (LR, DT, and RF) (Wu et al., 2022). Multiple parameters are
689 included in the models such as temperature, pH, oxidation-reduction potential, electrical

690 conductance, resistivity, total dissolved solids, salinity, pressure, DO, turbidity, and 24-h
691 accumulated rainfall. The results show that the RF model outperforms the other models in accuracy,
692 precision, recall, and AUC metrics.

693

694 **6.4.2. To track the sources of antimicrobial resistance genes in drinking water**

695 The challenges in detecting ARB and ARGs in drinking water with conventional methods highlight
696 the need for innovative solutions such as ML (Wu et al., 2022). ML effectively classifies the risks
697 of AMR and showcases the potential of AI technology in this area. Built on this technological
698 advancement, recent studies have employed another ML tool, SourceTracker (Knights et al., 2011),
699 which is based on Bayesian classification algorithm, to identify the potential sources of ARGs
700 (Chen et al., 2019; Wang et al., 2023). SourceTracker identified complex sources of ARGs and
701 assessed their contributions to ARG pollution in a peri-urban river (Chen et al., 2019). Results
702 show that the discharge from sewage treatment plants was the largest contributor of ARGs (81.6%
703 to 92.1%) in the river sediments. In another work, SourceTracker identified the presence of ARGs
704 in household drinking water and also traced their origins back to anthropogenic sources,
705 highlighting the significant impact of human activities on drinking water quality (Figure 5c) (Wang
706 et al., 2023). The data generated by SourceTracker have a strong Pearson correlation ($r = 0.98$)
707 with the corresponding expected proportion by artificial source inputs. Source tracking analysis
708 from that study indicates that a significant proportion of ARGs (37.1%) was from anthropogenic
709 sources, especially wastewater treatment plants.

710

711 **6.5. To study microbial communities in drinking water**

712 Both supervised and unsupervised ML algorithms are useful in metagenomic studies of microbial
713 communities in drinking water (Mahajna et al., 2022). When using the NB theorem for water
714 source tracking, one could apply either maximum posterior probability (Ritter et al., 2003)
715 (evaluation metrics: $RMSE_c$) or direct averaging posterior probability (Greenberg et al., 2010)
716 (evaluation metrics: $RMSE_p$) to estimate the distribution of microbes in water sources. Direct
717 averaging of the source posterior probability yields more precise source distribution estimates with
718 $RMSE_c$ being significantly lower than $RMSE_p$. The more precise source distribution estimates are
719 because direct estimation bypasses the information loss that typically happens when frequencies
720 are first classified and then averaged (Greenberg et al., 2010). SourceTracker is a ML tool
721 estimating the proportion of contaminants (Knights et al., 2011). That tool employs the Gibbs
722 sampling technique within a Bayesian framework and is more efficient than both the
723 aforementioned NB and RF classification-based source tracking methods (Smith et al., 2010). The
724 superior performance of SourceTracker is because it can handle ambiguity in the source and sink
725 distributions and can model a sink sample as a blend of various sources. SourceTracker can track
726 the origin of bacteria in tap water (Liu et al., 2018). For instance, a study created six ML models
727 to predict microbial contamination in a watershed using data on land cover, weather, and
728 hydrologic variables (Wu et al., 2020). In that case, SourceTracker generates ground-truth data for
729 training purposes. In addition, XGBoost outperforms the other five models (KNN, NB, SVM, NN,
730 and RF) in terms of accuracy and AUC when tracking the primary sources of microbial
731 contamination.

732

733 Unsupervised ML can unveil hidden features, trends, or patterns in bacterial communities in
734 engineered water systems. For instance, alpha and beta diversity analyses can display the spatial

735 dynamics and temporal trends of bacterial communities in DWDSs (Pinto et al., 2014). UniFrac
736 as another unsupervised ML tool uses a principal coordinates analysis (PCoA) coupled beta
737 diversity measure to analyze the differences among microbial communities in environmental
738 samples (Lozupone et al., 2011). UniFrac can effectively analyze the microbiome in drinking water
739 (Bruno et al., 2018; Li et al., 2017; Ling et al., 2018).

740

741 **7. Machine learning to detect accidental drinking water contamination**

742 **7.1. To detect anomalies and provide early warning of contamination in drinking water**

743 The increasing use of ML in safeguarding drinking water quality has led to the development of
744 innovative approaches to detect accidental drinking water contamination (Table 4) (Zhong et al.,
745 2021). This section explores the advancements of these ML approaches, underscoring the
746 versatility and potential of ML in ensuring safe drinking water supply by detecting accidental
747 drinking water contamination.

748

749 Three studies developed various ML approaches for anomaly detection on the basis of the drinking
750 water quality datasets from GECCO Industrial Challenges (GECCO IC) (Fehst et al., 2018;
751 Muharemi et al., 2019; Qian et al., 2020). These studies pinpointed shifts or spotted anomalies in
752 water quality over time. Various parameters such as pH, redox potential, electric conductivity,
753 turbidity, and chlorine dioxide concentration are the predictors, whereas events in Boolean form
754 are the outputs. A study developed SVM, DNN, long short-term memory (LSTM), recurrent NN
755 (RNN), LR, simple NN, and linear discriminant analysis (LDA) to detect water quality anomaly
756 in the dataset from 2017 GECCO IC (Muharemi et al., 2019). SVM shows the highest performance
757 with an F1-score of 0.99 in cross-validation. Nevertheless, all the models exhibit poor performance

758 with the unseen test dataset with a maximum F1-score of 0.36. In the other two studies focusing
759 on the dataset from GECCO IC 2018, LSTM demonstrates superior results, scoring a higher F1-
760 score than traditional models such as LR and SVM with F1-scores of 0.80 and 0.78, respectively
761 (Fehst et al., 2018; Qian et al., 2020).

762
763 The existing research, including the aforementioned key studies using the GECCO IC datasets,
764 has made significant strides in understanding anomaly detection in drinking water through ML. A
765 notable trend in recent research is the focus on real-time or online applications, reflecting a crucial
766 evolution toward practical, real-world implementations. Specifically, a study implemented an
767 LSTM-based approach to detect anomalies in water quality focusing on turbidity and conductivity
768 (Rodriguez-Perez et al., 2020) (Figure 6a). That study highlights the efficacy of semi-supervised
769 classification, which retains only normal values, in identifying abrupt changes and minor spikes in
770 water quality. By contrast, supervised classification, which considers both normal and anomalous
771 data, is more adept at identifying long-term anomalies linked to gradual changes. Notably, the
772 LSTM-based approach surpasses regression-based ARIMA in detecting these long-term anomalies.
773 Another study introduced an innovative stacking ensemble model designed for contamination
774 detection (Li et al., 2022). The model uses various water quality parameters such as total chlorine,
775 pH, electrical conductivity, temperature, TOC, and turbidity. That model integrates multiple ML
776 base predictors with a meta-predictor, trained through cross-validation. That approach enhances
777 the ability of the model to discern distinct features across water quality parameters. The ensemble
778 comprises base predictors such as ANN, SVM with a linear kernel, linear regressor, extra trees,
779 uniform weighted KNN, and a RF meta-predictor. The ensemble demonstrates superior
780 performance in detecting contamination compared with an ANN benchmark method, achieving

781 higher accuracy, lower false positive rates, and improved F1-scores. However, these models focus
782 on single-site, one-dimensional time series data, neglecting the spatial relationships inherent in
783 multi-site sensor data. This limitation could increase false alarm rates, particularly under
784 conditions of high hydraulic variability. To address this issue, a follow-up study proposed a novel
785 unsupervised, generative-adversarial-networks-based (GANs-based) multivariate multi-site
786 contamination event detection method (Figure 6b) (Li et al., 2023b). That method effectively
787 captures spatiotemporal patterns by transforming water quality data from single and multiple sites
788 into superimposed images. The GANs-based model, comprising a generator and a discriminator,
789 evaluates the degree of abnormality at each time step by generating anomaly scores. The generator
790 is trained to map historical image data to expected current images, while the discriminator
791 differentiates between generated and actual normal images. That method is benchmarked against
792 a multivariate unsupervised method using a minimum-volume-ellipsoid (MVE)-based event
793 detection model (Oliker and Ostfeld, 2014). That method demonstrates superior performance in
794 all contamination scenarios, including enhanced detection rates and reduced false alarms,
795 particularly for sensor groups positioned at varying distances from the contamination source.
796 Another unique ML approach can rapidly signal potential contamination risks (Asheri Arnon et al.,
797 2019). That approach uses an algorithm for the early detection of drinking water contamination
798 against an unpredictable stochastic background. By extracting key features from the
799 spectrophotometric characteristics of water, the algorithm can effectively identify contamination
800 using a unique affinity measure (Asheri-Arnon et al., 2018). The measure compares the absorbance
801 spectra of different water sources, thereby amplifying the feature dissimilarity between portable
802 and contaminated water, followed by processing via SVM and post-processing. That chain of data
803 processing generates a reliable early warning for contamination events with low false positives

804 and high true alarm accuracy. The pre-processing stage (the affinity measure and amplification) is
805 essential to achieving high accuracy but may be unnecessary for obtaining minimal false positives.

806

807 **8. Machine learning model distribution in safe drinking water supply**

808 We have prepared a macroscopic visual illustration to display the distribution of ML models across
809 research topics and sub-topics in safe drinking water supply (Figure 7). To facilitate a clear but
810 concise visual representation, we group certain ML models under broader principal categories on
811 the basis of their foundational architecture. For instance, models such as GA-ANN, Multi-layered-
812 ANN, and DNN share foundational characteristics inherent with ANN. Consequently, to elucidate
813 the overarching trends in model preferences across studies, we categorize these models as “NN-
814 based.” This approach discerns the broader trends and preferences in ML model selection and also
815 highlights the potential commonalities across research endeavors.

816

817 NN-based and regression-based models are the top two models frequently implemented. NN-based
818 models have significant applications in accessing and controlling DBPs and managing the
819 production and demand of drinking water. The prominent role of NN-based models in these two
820 fields is not coincidental but rather from the synergy between the inherent characteristics of these
821 fields and the strengths of NNs. Water management and DBP assessment often involve
822 multifaceted, nonlinear, and high-dimensional data sets that demand a robust modeling approach
823 (Aliashrafi et al., 2021; Ates et al., 2022; Ghobadi and Kang, 2023). Given their capability to model
824 complex non-linear relationships and handle various intricate data, NNs are an optimal solution in
825 these contexts. For instance, the unpredictability and variability in water demand patterns or the
826 multifarious factors influencing DBP formation both require a model that can discern patterns from

827 large, intricate datasets (Ahmadpour et al., 2023; Avni et al., 2015). Furthermore, the flexibility of
828 NNs in accommodating changing inputs makes them promising in assessing the dynamic nature
829 of drinking water systems. The wide applications of NNs in safe drinking water supply are due to
830 this harmonious fit between the challenges posed by these fields and the advantages of NNs.

831
832 By contrast, while regression-based models are widely applied in drinking water research, they
833 have suboptimal performances in certain contexts (Almheiri et al., 2021; Deng et al., 2021; Hong
834 et al., 2020; Legube et al., 2004; Lin et al., 2020; Rodriguez-Perez et al., 2020). This is not to
835 undermine the value of regression models. However, their linear or predefined non-linear
836 structures may limit their effectiveness, especially when compared with the adaptive and intricate
837 abilities of NNs.

838
839 The superior performance of the NN-based models is widely acknowledged (Goodfellow et al.,
840 2016). These general strengths become pertinent when the NN-based models are applied to the
841 complexities of water research. First, unlike regression models which are limited by their linear or
842 defined non-linear structures, NNs capture intricate, non-linear associations. Second, the mutable
843 architecture of NNs allows them to modify their framework during training, optimizing alignment
844 with the inherent data distribution. Lastly, given abundant data, NNs excel in discerning subtle
845 data patterns because of their proficiency in processing high-dimensional input attributes, whereas
846 regression models may have tendencies of underfitting. This proficiency of NNs is further
847 enhanced by the use of techniques such as grid search for hyperparameter optimization,
848 particularly crucial in fine-tuning the performance of NNs because of their complex architectures
849 and the numerous parameters required (Daniel et al., 2023; Rodriguez-Perez et al., 2020).

850
851 CNN-based models represent a specialized subclass of NN-based models adept at discerning
852 patterns in images or other forms of multi-dimensional data sets (LeCun et al., 1989; Lecun et al.,
853 1998). Therefore, we list the CNN-based models out of the broader NN category (Figure 7). The
854 practical implications of CNN-based models are evident in water research: They can interpret 2D
855 fluorescence spectra and predict the formation of DBPs during drinking water disinfection (Peleato,
856 2022), classify microbes using cell-level scattering images from drinking water (Luo et al., 2021;
857 Xu et al., 2020), and identify cells in 3D drinking water biofilm images (Jelli et al., 2023).

858
859 Other ensemble approaches are also widely applied in safe drinking water supply such as RF
860 (Breiman, 2001), XGB (Chen and Guestrin, 2016), boosted decision trees (BDT) (Friedman, 2001),
861 and stacking model (Wolpert, 1992). The core strength of these ensemble techniques is their ability
862 to amalgamate predictions from several models, aiming to boost accuracy and diminish overfitting
863 (Hastie et al., 2009). In water research where data can be noisy, varied, and sometimes sparse, such
864 strategies are invaluable. Several comparative studies have delved into the performance nuances
865 of different ensemble models. A recurring observation in these investigations is that the slight edge
866 XGB outperforms RF (Abdi and Mazloom, 2022; Park et al., 2020; Wu et al., 2020). Furthermore,
867 BRT outperforms RF (Bagriacik et al., 2018). Interestingly, while XGB has consistent prowess,
868 LightGBM, another gradient boosting framework, outperforms XGB (Abdi and Mazloom, 2022).
869 Therefore, as gradient boosting algorithms continue to evolve, newer iterations such as LightGBM
870 might offer even more refined performance. However, while ensemble methods offer certain
871 advantages, their efficacy is not universally dominant across scenarios. The best model is often
872 contingent upon the nature of the problem, the characteristics of the data, and the specific

873 objectives of the study. Ensemble models, with their ability to amalgamate insights from multiple
874 “weak learners,” might excel in scenarios where data are diverse, noisy, or sparse (Fasel et al.,
875 2022; Pang et al., 2018; Sluban and Lavrač, 2015). By contrast, for problems where the data
876 structures are deeply hierarchical or when data patterns are straightforward, NN-based models or
877 regression models are more suitable. The crucial factor is to match the ability of the models with
878 the specific demands and characteristics of the data sets.

879

880 We include (S)ARIMA, Kriging interpolation, SaTScan, LDFA, Alpha and Beta diversity analyses,
881 UniFrac, and MVE in statistical models (Figure 7). These models are more deterministic and often
882 rooted in foundational principles and known relationships. For instance, SARIMA and Kriging
883 interpolation can capture temporal and spatial patterns, respectively (Guo et al., 2018; Tian et al.,
884 2020). Alpha and Beta diversity analyses and UniFrac quantify microbial community diversity and
885 compositional differences (Bruno et al., 2018; Li et al., 2017; Ling et al., 2018; Lozupone et al.,
886 2011; Pinto et al., 2014). These models typically operate under specific assumptions about the
887 underlying data distribution or spatiotemporal relationships. By contrast, ML, especially deep
888 learning, is more adaptive, learning patterns directly from the data without stringent assumptions
889 (Khattak et al., 2022; Savadatti et al., 2022; Singh et al., 2023).

890

891 **9. Challenges and outlooks for the applications of machine learning in safe drinking water** 892 **supply**

893 While ML has made significant advancements in drinking water research, several areas remain
894 untapped, offering significant potential for exploration and improvement. Crucial topics, such as
895 biofilm development, the assessment of AMR risks, and the evaluation of pathogen-related dangers

896 in engineered water systems, are not fully represented in current research. The untapped potential
897 in these fields is immense, and the need to bridge the interdisciplinary divide is critical.

898

899 One significant barrier is the disconnect between water treatment experts and AI specialists. Water
900 scientists and engineers may not be conversant with the nuances of AI, while AI technologists
901 might lack knowledge of water treatment and supply. This knowledge gap impedes the effective
902 deployment of ML in enhancing the safety of drinking water supply. Addressing this dichotomy is
903 beneficial and essential, necessitating educational and collaborative efforts to build a shared
904 understanding and to develop interdisciplinary skill sets.

905

906 Further complicating the matter is the absence of standardized tool kits tailored to the drinking
907 water sector. Such standardization is vital for enabling consistent application across various
908 research and implementation efforts. Uniformity in tools and approaches would not only
909 streamline the processes but also bolster collaborative work, which is often fragmented across
910 regions and specializations.

911

912 Advancements in ML tools must cater to the unique challenges presented by safe drinking water
913 supply. Water quality in engineered water systems is affected by numerous factors that vary
914 spatiotemporally, requiring ML solutions that can adapt to and learn from these dynamic conditions.
915 Thus, future studies should customize existing ML frameworks or innovate new ones that can
916 grapple with the complexities inherent to the water supply networks.

917

918 Looking to the horizon, the broader vision involves leveraging ML to address the global drinking
919 water crisis. Issues such as water scarcity, the presence of emerging contaminants, and the
920 formation of DBPs present a global challenge. ML tools have been predominantly developed with
921 local or regional contexts, yet the drinking water crisis demands a global perspective. The ambition
922 to harness ML for these universal challenges is not only visionary but also an imperative stride
923 toward ensuring water security worldwide.

924
925 In pursuit of these goals, the integration of advanced ML models becomes a cornerstone in tackling
926 the multifaceted issues tied to drinking water safety. Future endeavors should also prioritize the
927 promotion of open-access data sharing within and beyond the drinking water research community
928 (Zhong et al., 2021). The endeavors will enhance collaboration, drive transparency, and support
929 the reproducibility of scientific findings, which are the bedrock of robust research. Furthermore,
930 establishing a comprehensive comparative framework to evaluate different ML models will be
931 instrumental in identifying the optimal solutions for the challenges in drinking water research. By
932 embracing these strategies, we can aspire to not just bridge existing knowledge gaps but also
933 significantly elevate the role of ML in securing safe and more sustainable water supply.

934

935 **9. Conclusions**

936 Assessing and ensuring safe drinking water supply is a global challenge with conventional
937 approaches. ML as a novel methodology is promising in monitoring and protecting drinking water
938 quality, especially in municipal engineered water systems. This review for the first time
939 comprehensively summarizes the applications of ML in assessing and ensuring safe drinking water
940 supply with a focus on water quality in engineered water systems. We compile the applications of

941 ML from the physical, chemical, and microbiological perspectives. From the physical perspective,
942 ML is useful in managing drinking water production and demand and monitoring drinking water
943 pipeline failures. From the chemical perspective, ML is promising in assessing and controlling
944 DBPs, monitoring and mitigating heavy metals, and tracking nitrification in drinking water. From
945 the microbiological perspective, ML can monitor and mitigate OPs, detect *Cryptosporidium* and
946 *Giardia*, assess biofilm development, predict AMR risks, and study microbial communities in
947 municipal water, especially in engineered water systems. In addition, ML is a useful tool in
948 detecting accidental drinking water contamination.

949

950 **Acknowledgment**

951 This work has been supported by the Faculty Research Fund from Arkansas State University
952 (Jonesboro, AR, U.S.A.).

953

954 **Declaration of interest**

955 No potential or actual conflict of interest exists in this work.

956

957

958

959 **Uncategorized References**

960 Abdi, J. and Mazloom, G. 2022. Machine learning approaches for predicting arsenic adsorption
961 from water using porous metal-organic frameworks. *Sci Rep* 12(1), 16458.

962 Adeyemo, F.E., Singh, G., Reddy, P., Bux, F. and Stenstrom, T.A. 2019. Efficiency of chlorine
963 and UV in the inactivation of *Cryptosporidium* and *Giardia* in wastewater. *PLoS One* 14(5),
964 e0216040.

965 Ahmadpour, E., Delpla, I., Debia, M., Simard, S., Proulx, F., Sérodes, J.-B., Valois, I., Tardif, R.,
966 Haddad, S. and Rodriguez, M. 2023. Full-scale multisampling and empirical modeling of

- 967 DBPs in water and air of indoor pools. *Environmental Monitoring and Assessment* 195(9),
968 1128.
- 969 Ahmed, A.N., Othman, F.B., Afan, H.A., Ibrahim, R.K., Fai, C.M., Hossain, M.S., Ehteram, M.
970 and Elshafie, A. 2019. Machine learning methods for better water quality prediction. *Journal*
971 *of Hydrology* 578, 124084.
- 972 Aliashrafi, A., Zhang, Y., Groenewegen, H. and Peleato, N.M. 2021. A review of data-driven
973 modelling in drinking water treatment. *Reviews in Environmental Science and Bio/Technology*
974 20(4), 985-1009.
- 975 Allesoe, R.L., Lundgaard, A.T., Hernandez Medina, R., Aguayo-Orozco, A., Johansen, J., Nissen,
976 J.N., Brorsson, C., Mazzoni, G., Niu, L., Biel, J.H., Leal Rodriguez, C., Brasas, V., Webel, H.,
977 Benros, M.E., Pedersen, A.G., Chmura, P.J., Jacobsen, U.P., Mari, A., Koivula, R., Mahajan,
978 A., Vinuela, A., Tajes, J.F., Sharma, S., Haid, M., Hong, M.G., Musholt, P.B., De Masi, F., Vogt,
979 J., Pedersen, H.K., Gudmundsdottir, V., Jones, A., Kennedy, G., Bell, J., Thomas, E.L., Frost,
980 G., Thomsen, H., Hansen, E., Hansen, T.H., Vestergaard, H., Muilwijk, M., Blom, M.T., t Hart,
981 L.M., Pattou, F., Raverdy, V., Brage, S., Kokkola, T., Heggie, A., McEvoy, D., Mourby, M.,
982 Kaye, J., Hattersley, A., McDonald, T., Ridderstrale, M., Walker, M., Forgie, I., Giordano, G.N.,
983 Pavo, I., Ruetten, H., Pedersen, O., Hansen, T., Dermitzakis, E., Franks, P.W., Schwenk, J.M.,
984 Adamski, J., McCarthy, M.I., Pearson, E., Banasik, K., Rasmussen, S., Brunak, S. and
985 Consortium, I.D. 2023. Discovery of drug-omics associations in type 2 diabetes with
986 generative deep-learning models. *Nat Biotechnol* 41(3), 399-408.
- 987 Almalioglu, Y., Turan, M., Trigoni, N. and Markham, A. 2022. Deep learning-based robust
988 positioning for all-weather autonomous driving. *Nat Mach Intell* 4(9), 749-760.
- 989 Almheiri, Z., Meguid, M. and Zayed, T. 2021. Failure modeling of water distribution pipelines
990 using meta-learning algorithms. *Water Res* 205, 117680.
- 991 Asheri Arnon, T., Ezra, S. and Fishbain, B. 2019. Water characterization and early contamination
992 detection in highly varying stochastic background water, based on Machine Learning
993 methodology for processing real-time UV-Spectrophotometry. *Water Research* 155, 333-342.
- 994 Asheri-Arnon, T., Ezra, S. and Fishbain, B. 2018. Contamination Detection of Water with Varying
995 Routine Backgrounds by UV-Spectrophotometry. *Journal of Water Resources Planning and*
996 *Management* 144(9), 04018056.
- 997 Ates, N., Civelekoglu, G. and Kaplan-Bekaroglu, S.S. (2022) *Water and Wastewater Management:*
998 *Global Problems and Measures.* Bahadir, M. and Haarstrick, A. (eds), pp. 67-82, Springer
999 International Publishing, Cham.
- 1000 Avni, N., Fishbain, B. and Shamir, U. 2015. Water consumption patterns as a basis for water
1001 demand modeling. *Water Resources Research* 51(10), 8165-8181.
- 1002 Bagriacik, A., Davidson, R.A., Hughes, M.W., Bradley, B.A. and Cubrinovski, M. 2018.
1003 Comparison of statistical and machine learning approaches to modeling earthquake damage to
1004 water pipelines. *Soil Dynamics and Earthquake Engineering* 112, 76-88.
- 1005 Bakker, M., Vreeburg, J.H.G., Van De Roer, M. and Rietveld, L.C. 2014. Heuristic burst detection
1006 method using flow and pressure measurements. *Journal of Hydroinformatics* 16(5), 1194-1209.
- 1007 Balogun, A.-L., Tella, A., Baloo, L. and Adebisi, N. 2021. A review of the inter-correlation of
1008 climate change, air pollution and urban sustainability using novel machine learning algorithms
1009 and spatial information science. *Urban Climate* 40, 100989.
- 1010 Benedict, K.M., Reses, H., Vigar, M., Roth, D.M., Roberts, V.A., Mattioli, M., Cooley, L.A.,
1011 Hilborn, E.D., Wade, T.J. and Fullerton, K.E. 2017. Surveillance for waterborne disease

- 1012 outbreaks associated with drinking water-United States, 2013-2014. *Morbidity and Mortality*
1013 *Weekly Report* 66(44), 1216.
- 1014 Benítez, J.S., Rodríguez, C.M. and Casas, A.F. 2021. Disinfection byproducts (DBPs) in drinking
1015 water supply systems: a systematic review. *Physics and Chemistry of the Earth, Parts A/B/C*
1016 123, 102987.
- 1017 Berglund, E., Vizanko, B., Kadinski, L. and Ostfeld, A. (2023) Coupling Machine Learning and
1018 Agent-Based Modeling to Characterize Contamination Sources in Water Distribution Systems
1019 for Changing Demand Regimes. *World Environmental and Water Resources Congress 2023*,
1020 pp. 881-890.
- 1021 Berne, C., Ellison, C.K., Ducret, A. and Brun, Y.V. 2018. Bacterial adhesion at the single-cell
1022 level. *Nature Reviews Microbiology* 16(10), 616-627.
- 1023 Bhagat, S.K., Tung, T.M. and Yaseen, Z.M. 2020. Development of artificial intelligence for
1024 modeling wastewater heavy metal removal: State of the art, application assessment and
1025 possible future research. *Journal of Cleaner Production* 250, 119473.
- 1026 Bi, K., Xie, L., Zhang, H., Chen, X., Gu, X. and Tian, Q. 2023. Accurate medium-range global
1027 weather forecasting with 3D neural networks. *Nature* 619(7970), 533-538.
- 1028 Bond, T., Goslan, E.H., Parsons, S.A. and Jefferson, B. 2011. Treatment of disinfection by-
1029 product precursors. *Environmental Technology* 32(1), 1-25.
- 1030 Bond, T. and Graham, N. 2017. Predicting chloroform production from organic precursors. *Water*
1031 *Research* 124, 167-176.
- 1032 Bond, T., Templeton, M.R. and Graham, N. 2012. Precursors of nitrogenous disinfection by-
1033 products in drinking water—A critical review and analysis. *Journal of Hazardous Materials*
1034 235-236, 1-16.
- 1035 Bradley, A.P. 1997. The use of the area under the ROC curve in the evaluation of machine learning
1036 algorithms. *Pattern Recognition* 30(7), 1145-1159.
- 1037 Breiman, L. 2001. Random Forests. *Machine Learning* 45(1), 5-32.
- 1038 Brunello, A., Civilini, M., De Martin, S., Felice, A., Franchi, M., Iacumin, L., Saccomanno, N. and
1039 Vitacolonna, N. 2022. Machine learning-assisted environmental surveillance of *Legionella*:
1040 A retrospective observational study in Friuli-Venezia Giulia region of Italy in the period 2002–
1041 2019. *Informatics in Medicine Unlocked* 28, 100803.
- 1042 Bruno, A., Sandionigi, A., Bernasconi, M., Panio, A., Labra, M. and Casiraghi, M. 2018. Changes
1043 in the Drinking Water Microbiome: Effects of Water Treatments Along the Flow of Two
1044 Drinking Water Treatment Plants in a Urbanized Area, Milan (Italy). *Front Microbiol* 9, 2557.
- 1045 Cabral, J.P. 2010. Water microbiology. Bacterial pathogens and water. *International Journal of*
1046 *Environmental Research and Public Health* 7(10), 3657-3703.
- 1047 CDC 2021a. Parasites - Cryptosporidium (also known as "Crypto").
- 1048 CDC 2021b. Parasites - Giardia.
- 1049 Cha, D., Park, S., Kim, M.S., Kim, T., Hong, S.W., Cho, K.H. and Lee, C. 2021. Prediction of
1050 Oxidant Exposures and Micropollutant Abatement during Ozonation Using a Machine
1051 Learning Method. *Environ Sci Technol* 55(1), 709-718.
- 1052 Chang, Q., Wang, W., Regev-Yochay, G., Lipsitch, M. and Hanage, W.P. 2015. Antibiotics in
1053 agriculture and the risk to human health: how worried should we be? *Evol Appl* 8(3), 240-247.
- 1054 Chen, H., Bai, X., Li, Y., Jing, L., Chen, R. and Teng, Y. 2019. Source identification of antibiotic
1055 resistance genes in a peri-urban river using novel crAssphage marker genes and metagenomic
1056 signatures. *Water Research* 167, 115098.

- 1057 Chen, T. and Guestrin, C. 2016 XGBoost: A Scalable Tree Boosting System, pp. 785–794,
1058 Association for Computing Machinery, San Francisco, California, USA.
- 1059 Chowdhury, S., Mazumder, M.A.J., Al-Attas, O. and Husain, T. 2016. Heavy metals in drinking
1060 water: Occurrences, implications, and future needs in developing countries. *Sci Total Environ*
1061 569-570, 476-488.
- 1062 Cordero, J.A., He, K., Janya, K., Echigo, S. and Itoh, S. 2021. Predicting formation of haloacetic
1063 acids by chlorination of organic compounds using machine-learning-assisted quantitative
1064 structure-activity relationships. *Journal of Hazardous Materials* 408, 124466.
- 1065 Craun, G.F., Brunkard, J.M., Yoder, J.S., Roberts, V.A., Carpenter, J., Wade, T., Calderon, R.L.,
1066 Roberts, J.M., Beach, M.J. and Roy, S.L. 2010. Causes of outbreaks associated with drinking
1067 water in the United States from 1971 to 2006. *Clinical Microbiology Reviews* 23(3), 507-528.
- 1068 Cutler, D. and Miller, G. 2005. The role of public health improvements in health advances: the
1069 twentieth-century United States. *Demography* 42, 1-22.
- 1070 Daniel, I., Abhijith, G.R., Kadinski, L., Ostfeld, A. and Cominola, A. (2023) A Machine Learning-
1071 Based Surrogate Model for Coupled Hydraulic and Water Quality Simulation in Water
1072 Distribution Networks. *World Environmental and Water Resources Congress 2023*, pp. 817-
1073 830.
- 1074 De Jesus, K.L.M., Senoro, D.B., Dela Cruz, J.C. and Chan, E.B. 2021. A Hybrid Neural Network-
1075 Particle Swarm Optimization Informed Spatial Interpolation Technique for Groundwater
1076 Quality Mapping in a Small Island Province of the Philippines. *Toxics* 9(11).
- 1077 Delpla, I., Jung, A.-V., Baures, E., Clement, M. and Thomas, O. 2009. Impacts of climate change
1078 on surface water quality in relation to drinking water production. *Environment International*
1079 35(8), 1225-1233.
- 1080 Deng, Y., Zhou, X., Shen, J., Xiao, G., Hong, H., Lin, H., Wu, F. and Liao, B.-Q. 2021. New
1081 methods based on back propagation (BP) and radial basis function (RBF) artificial neural
1082 networks (ANNs) for predicting the occurrence of haloketones in tap water. *Science of The*
1083 *Total Environment* 772, 145534.
- 1084 Dogo, E.M., Nwulu, N.I., Twala, B. and Aigbavboa, C. 2019. A survey of machine learning
1085 methods applied to anomaly detection on drinking-water quality data. *Urban Water Journal*
1086 16(3), 235-248.
- 1087 Donohue, M., King, D., Pfaller, S. and Mistry, J. 2019. The sporadic nature of *Legionella*
1088 *pneumophila*, *Legionella pneumophila* Sg1 and *Mycobacterium avium* occurrence within
1089 residences and office buildings across 36 states in the United States. *Journal of applied*
1090 *microbiology* 126(5), 1568-1579.
- 1091 Erhan, D., Szegedy, C., Toshev, A. and Anguelov, D. 2014 Scalable Object Detection Using Deep
1092 Neural Networks, pp. 2155–2162, IEEE Computer Society.
- 1093 Ewuzie, U., Bolade, O.P. and Egbedina, A.O. (2022) Current Trends and Advances in Computer-
1094 Aided Intelligent Environmental Data Engineering. Marques, G. and Ighalo, J.O. (eds), pp.
1095 185-218, Elsevier Inc., San Diego, California, U.S.A.
- 1096 Fan, X., Zhang, X., Yu, A., Speitel, M. and Yu, X. 2023. Assessment of the impacts of climat
1097 change on water supply system pipe failures. *Scientific Reports* 13(1), 7349.
- 1098 Fasel, U., Kutz, J.N., Brunton, B.W. and Brunton, S.L. 2022. Ensemble-SINDy: Robust sparse
1099 model discovery in the low-data, high-noise limit, with active learning and control.
1100 *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*
1101 478(2260), 20210904.

- 1102 Favere, J., Barbosa, R.G., Sleutels, T., Verstraete, W., De Gusseme, B. and Boon, N. 2021.
1103 Safeguarding the microbial water quality from source to tap. *Npj Clean Water* 4, 28.
- 1104 Fehst, V., La, H.C., Nghiem, T.-D., Mayer, B.E., Englert, P. and Fiebig, K.-H. 2018 Automatic vs.
1105 manual feature engineering for anomaly detection of drinking-water quality, pp. 5–6,
1106 Association for Computing Machinery, Kyoto, Japan.
- 1107 Feng, S., Sun, H., Yan, X., Zhu, H., Zou, Z., Shen, S. and Liu, H.X. 2023. Dense reinforcement
1108 learning for safety validation of autonomous vehicles. *Nature* 615(7953), 620-627.
- 1109 Folkman, S. 2018. Water Main Break Rates in the USA and Canada: A Comprehensive Study,
1110 2018. Mechanical and Aerospace Engineering Faculty Publications Paper 174.
- 1111 Friedman, J.H. 2001. Greedy function approximation: A gradient boosting machine. *The Annals*
1112 *of Statistics* 29(5), 1189-1232, 1144.
- 1113 Ghiassi, M., Fa'al, F. and Abrishamchi, A. 2017. Large metropolitan water demand forecasting
1114 using DAN2, FTDNN, and KNN models: A case study of the city of Tehran, Iran. *Urban Water*
1115 *Journal* 14(6), 655-659.
- 1116 Ghiassi, M., Saidane, H. and Zimbra, D.K. 2005. A dynamic artificial neural network model for
1117 forecasting time series events. *International Journal of Forecasting* 21(2), 341-362.
- 1118 Ghobadi, F. and Kang, D. 2023. Application of Machine Learning in Water Resources
1119 Management: A Systematic Literature Review. *Water* 15(4), 620.
- 1120 Gogoi, A., Mazumder, P., Tyagi, V.K., Chaminda, G.T., An, A.K. and Kumar, M. 2018.
1121 Occurrence and fate of emerging contaminants in water environment: a review. *Groundwater*
1122 *for Sustainable Development* 6, 169-180.
- 1123 Goh, S.G., Jiang, P., Ng, C., Le, T.H., Haller, L., Chen, H., Charles, F.R., Chen, H., Liu, X., He, Y.
1124 and Gin, K.Y. 2022. A new modelling framework for assessing the relative burden of
1125 antimicrobial resistance in aquatic environments. *J Hazard Mater* 424(Pt C), 127621.
- 1126 Gomez-Alvarez, V. and Revetta, R.P. 2020. Monitoring of Nitrification in Chloraminated
1127 Drinking Water Distribution Systems With Microbiome Bioindicators Using Supervised
1128 Machine Learning. *Front Microbiol* 11, 571009.
- 1129 Gong, J., Guo, X., Yan, X. and Hu, C. 2023. Review of urban drinking water contamination
1130 source identification methods. *Energies* 16(2), 705.
- 1131 Gonzalez, S., Lopez-Roldan, R. and Cortina, J.-L. 2013. Presence of metals in drinking water
1132 distribution networks due to pipe material leaching: a review. *Toxicological & Environmental*
1133 *Chemistry* 95(6), 870-889.
- 1134 Goodfellow, I., Bengio, Y. and Courville, A. (2016) *Deep Learning*, MIT Press.
- 1135 Greenberg, J., Price, B. and Ware, A. 2010. Alternative estimate of source distribution in microbial
1136 source tracking using posterior probabilities. *Water Res* 44(8), 2629-2637.
- 1137 Grey, D. and Sadoff, C.W. 2007. Sink or Swim? Water security for growth and development.
1138 *Water Policy* 9(6), 545-571.
- 1139 Guo, G., Liu, S., Wu, Y., Li, J., Zhou, R. and Zhu, X. 2018. Short-Term Water Demand Forecast
1140 Based on Deep Learning Method. *Journal of Water Resources Planning and Management*
1141 144(12), 04018076.
- 1142 Gupta, S. and Basant, N. 2016. Modeling the reactivity of ozone and sulphate radicals towards
1143 organic chemicals in water using machine learning approaches. *RSC Advances* 6(110),
1144 108448-108457.
- 1145 Gupta, V.K. and Ali, I. (2013) *Environmental Water*. Gupta, V.K. and Ali, I. (eds), pp. 1-27, Elsevier.

- 1146 Hamidian, A.H., Esfandeh, S., Zhang, Y. and Yang, M. 2019. Simulation and optimization of
1147 nanomaterials application for heavy metal removal from aqueous solutions. *Inorganic and*
1148 *Nano-Metal Chemistry* 49(7), 217-230.
- 1149 Hastie, T., Tibshirani, R. and Friedman, J.H. (2009) *The elements of statistical learning : data*
1150 *mining, inference, and prediction*, Springer, New York.
- 1151 He, L., Bai, L., Dionysiou, D.D., Wei, Z., Spinney, R., Chu, C., Lin, Z. and Xiao, R. 2021.
1152 Applications of computational chemistry, artificial intelligence, and machine learning in
1153 aquatic chemistry research. *Chemical Engineering Journal* 426, 131810.
- 1154 Helte, E., Säve-Söderbergh, M., Larsson, S.C., Martling, A. and Åkesson, A. 2023. Disinfection
1155 by-products in drinking water and risk of colorectal cancer: a population-based cohort study.
1156 *Journal of the National Cancer Institute*, djad145.
- 1157 Henrique Alves Ribeiro, V. and Reynoso-Meza, G. 2023. Multi-criteria Decision-Making
1158 Techniques for the Selection of Pareto-optimal Machine Learning Models in a Drinking-Water
1159 Quality Monitoring Problem. *International Journal of Information Technology & Decision*
1160 *Making*, 1-28.
- 1161 Hong, H., Zhang, Z., Guo, A., Shen, L., Sun, H., Liang, Y., Wu, F. and Lin, H. 2020. Radial basis
1162 function artificial neural network (RBF ANN) as well as the hybrid method of RBF ANN and
1163 grey relational analysis able to well predict trihalomethanes levels in tap water. *Journal of*
1164 *Hydrology* 591, 125574.
- 1165 Hossain, S., Cook, D., Chow, C.W.K. and Hewa, G.A. 2021. Development of an Optical Method
1166 to Monitor Nitrification in Drinking Water. *Sensors (Basel)* 21(22).
- 1167 Hu, G., Mian, H.R., Mohammadiun, S., Rodriguez, M.J., Hewage, K. and Sadiq, R. 2023.
1168 Appraisal of machine learning techniques for predicting emerging disinfection byproducts in
1169 small water distribution networks. *Journal of Hazardous Materials* 446, 130633.
- 1170 Huang, R., Ma, C., Ma, J., Huangfu, X. and He, Q. 2021. Machine learning in natural and
1171 engineered water systems. *Water Res* 205, 117666.
- 1172 Huang, Y., Li, T., Zheng, S., Fan, L., Su, L., Zhao, Y., Xie, H.B. and Li, C. 2020. QSAR modeling
1173 for the ozonation of diverse organic compounds in water. *Sci Total Environ* 715, 136816.
- 1174 Hutchings, M.I., Truman, A.W. and Wilkinson, B. 2019. Antibiotics: past, present and future.
1175 *Current Opinion in Microbiology* 51, 72-80.
- 1176 Isaac, T.S. and Sherchan, S.P. 2020. Molecular detection of opportunistic premise plumbing
1177 pathogens in rural Louisiana's drinking water distribution system. *Environmental Research* 181,
1178 108847.
- 1179 James, G., Witten, D., Hastie, T. and Tibshirani, R. 2013 *An introduction to statistical learning*
1180 *with applications in R*, Springer,, New York.
- 1181 Jefri, U.H.N.M., Khan, A., Lim, Y.C., Lee, K.S., Liew, K.B., Kassab, Y.W., Choo, C.-Y., Al-Worafi,
1182 Y.M., Ming, L.C. and Kalusalingam, A. 2022. A systematic review on chlorine dioxide as a
1183 disinfectant. *Journal of Medicine and Life* 15(3), 313.
- 1184 Jelli, E., Ohmura, T., Netter, N., Abt, M., Jiménez-Siebert, E., Neuhaus, K., Rode, D.K.H., Nadell,
1185 C.D. and Drescher, K. 2023. Single-cell segmentation in bacterial biofilms with an optimized
1186 deep learning method enables tracking of cell lineages and measurements of growth rates.
1187 *Molecular Microbiology* 119(6), 659-676.
- 1188 Jia, X., O'Connor, D., Shi, Z. and Hou, D. 2021. VIRS based detection in combination with
1189 machine learning for mapping soil pollution. *Environ Pollut* 268(Pt A), 115845.
- 1190 Jiang, S., Hu, J., Wood, K.L. and Luo, J. 2021. Data-Driven Design-By-Analogy: State-of-the-
1191 Art and Future Directions. *Journal of Mechanical Design* 144(2).

- 1192 Joseph, L., Jun, B.-M., Flora, J.R.V., Park, C.M. and Yoon, Y. 2019. Removal of heavy metals
1193 from water sources in the developing world using low-cost materials: A review. *Chemosphere*
1194 229, 142-159.
- 1195 Kaufmann, E., Bauersfeld, L., Loquercio, A., Muller, M., Koltun, V. and Scaramuzza, D. 2023.
1196 Champion-level drone racing using deep reinforcement learning. *Nature* 620(7976), 982-987.
- 1197 Kazemi, E., Kyritsakas, G., Husband, S., Flavell, K., Speight, V. and Boxall, J. 2023. Predicting
1198 iron exceedance risk in drinking water distribution systems using machine learning. *IOP*
1199 *Conference Series: Earth and Environmental Science* 1136(1), 012047.
- 1200 Khattak, A., Bukhsh, R., Aslam, S., Yafoz, A., Alghushairy, O. and Alsini, R. 2022. A Hybrid
1201 Deep Learning-Based Model for Detection of Electricity Losses Using Big Data in Power
1202 Systems. *Sustainability* 14(20), 13627.
- 1203 Kirstein, I.V., Hensel, F., Gomiero, A., Iordachescu, L., Vianello, A., Wittgren, H.B. and Vollertsen,
1204 J. 2021. Drinking plastics?—Quantification and qualification of microplastics in drinking
1205 water distribution systems by μ FTIR and Py-GCMS. *Water Research* 188, 116519.
- 1206 Knights, D., Kuczynski, J., Charlson, E.S., Zaneveld, J., Mozer, M.C., Collman, R.G., Bushman,
1207 F.D., Knight, R. and Kelley, S.T. 2011. Bayesian community-wide culture-independent
1208 microbial source tracking. *Nat Methods* 8(9), 761-763.
- 1209 Krasner, S.W., Mitch, W.A., McCurry, D.L., Hanigan, D. and Westerhoff, P. 2013. Formation,
1210 precursors, control, and occurrence of nitrosamines in drinking water: a review. *Water Res*
1211 47(13), 4433-4450.
- 1212 Krasner, S.W., Weinberg, H.S., Richardson, S.D., Pastor, S.J., Chinn, R., Scilimenti, M.J., Onstad,
1213 G.D. and Thruston, A.D. 2006. Occurrence of a New Generation of Disinfection Byproducts.
1214 *Environmental Science & Technology* 40(23), 7175-7185.
- 1215 LeCun, Y., Boser, B., Denker, J.S., Henderson, D., Howard, R.E., Hubbard, W. and Jackel, L.D.
1216 1989. Backpropagation Applied to Handwritten Zip Code Recognition. *Neural Computation*
1217 1(4), 541-551.
- 1218 Lecun, Y., Bottou, L., Bengio, Y. and Haffner, P. 1998. Gradient-based learning applied to
1219 document recognition. *Proceedings of the IEEE* 86(11), 2278-2324.
- 1220 Lee, D., Gibson, J.M., Brown, J., Habtewold, J. and Murphy, H.M. 2023. Burden of disease from
1221 contaminated drinking water in countries with high access to safely managed water: a
1222 systematic review. *Water Research* 242, 120244.
- 1223 Lee, J., Kim, E.-S., Roh, B.-S., Eom, S.-W. and Zoh, K.-D. 2013. Occurrence of disinfection by-
1224 products in tap water distribution systems and their associated health risk. *Environmental*
1225 *Monitoring and Assessment* 185, 7675-7691.
- 1226 Legube, B., Parinet, B., Gelinet, K., Berne, F. and Croue, J.P. 2004. Modeling of bromate
1227 formation by ozonation of surface waters in drinking water treatment. *Water Res* 38(8), 2185-
1228 2195.
- 1229 Leitão, J., Simões, N., Sá Marques, J.A., Gil, P., Ribeiro, B. and Cardoso, A. 2019. Detecting
1230 urban water consumption patterns: a time-series clustering approach. *Water Supply* 19(8),
1231 2323-2329.
- 1232 Ley, C., Martin, R.K., Pareek, A., Groll, A., Seil, R. and Tischer, T. 2022. Machine learning and
1233 conventional statistics: making sense of the differences. *Knee Surg Sports Traumatol Arthrosc*
1234 30(3), 753-757.
- 1235 Li, J., Heap, A.D., Potter, A. and Daniell, J.J. 2011. Application of machine learning methods to
1236 spatial interpolation of environmental variables. *Environmental Modelling & Software* 26(12),
1237 1647-1659.

- 1238 Li, J., Zhang, Z., Xiang, Y., Jiang, J. and Yin, R. 2023a. Role of UV-based advanced oxidation
1239 processes on NOM alteration and DBP formation in drinking water treatment: A state-of-the-
1240 art review. *Chemosphere* 311, 136870.
- 1241 Li, L., Rong, S., Wang, R. and Yu, S. 2021. Recent advances in artificial intelligence and machine
1242 learning for nonlinear relationship analysis and process control in drinking water treatment: A
1243 review. *Chemical Engineering Journal* 405, 126673.
- 1244 Li, P. and Wu, J. 2019. Drinking water quality and public health. *Exposure and Health* 11(2), 73-
1245 79.
- 1246 Li, Q., Yu, S., Li, L., Liu, G., Gu, Z., Liu, M., Liu, Z., Ye, Y., Xia, Q. and Ren, L. 2017. Microbial
1247 Communities Shaped by Treatment Processes in a Drinking Water Treatment Plant and Their
1248 Contribution and Threat to Drinking Water Safety. *Front Microbiol* 8, 2465.
- 1249 Li, R.A., McDonald, J.A., Sathasivan, A. and Khan, S.J. 2019. Disinfectant residual stability
1250 leading to disinfectant decay and by-product formation in drinking water distribution systems:
1251 A systematic review. *Water Res* 153, 335-348.
- 1252 Li, Z., Liu, H., Zhang, C. and Fu, G. 2023b. Generative adversarial networks for detecting
1253 contamination events in water distribution systems using multi-parameter, multi-site water
1254 quality monitoring. *Environmental Science and Ecotechnology* 14, 100231.
- 1255 Li, Z., Zhang, C., Liu, H., Zhang, C., Zhao, M., Gong, Q. and Fu, G. 2022. Developing stacking
1256 ensemble models for multivariate contamination detection in water distribution systems.
1257 *Science of The Total Environment* 828, 154284.
- 1258 Ligda, P., Claerebout, E., Kostopoulou, D., Zdragas, A., Casaert, S., Robertson, L.J. and Sotiraki,
1259 S. 2020. Cryptosporidium and Giardia in surface water and drinking water: Animal sources
1260 and towards the use of a machine-learning approach as a tool for predicting contamination.
1261 *Environ Pollut* 264, 114766.
- 1262 Lin, H., Dai, Q., Zheng, L., Hong, H., Deng, W. and Wu, F. 2020. Radial basis function artificial
1263 neural network able to accurately predict disinfection by-product levels in tap water: Taking
1264 haloacetic acids as a case study. *Chemosphere* 248, 125999.
- 1265 Lin, T.Y., Dollár, P., Girshick, R., He, K., Hariharan, B. and Belongie, S. 2017. Feature Pyramid
1266 Networks for Object Detection, pp. 936-944.
- 1267 Ling, F., Whitaker, R., LeChevallier, M.W. and Liu, W.T. 2018. Drinking water microbiome
1268 assembly induced by water stagnation. *ISME J* 12(6), 1520-1531.
- 1269 Liu, G., Zhang, Y., van der Mark, E., Magic-Knezev, A., Pinto, A., van den Bogert, B., Liu, W.,
1270 van der Meer, W. and Medema, G. 2018. Assessing the origin of bacteria in tap water and
1271 distribution system in an unchlorinated drinking water system by SourceTracker using
1272 microbial community fingerprints. *Water Res* 138, 86-96.
- 1273 Liu, S., Gunawan, C., Barraud, N., Rice, S.A., Harry, E.J. and Amal, R. 2016. Understanding,
1274 Monitoring, and Controlling Biofilm Growth in Drinking Water Distribution Systems. *Environ*
1275 *Sci Technol* 50(17), 8954-8976.
- 1276 Lowe, M., Qin, R. and Mao, X. 2022. A review on machine learning, artificial intelligence, and
1277 smart technology in water treatment and monitoring. *Water* 14(9), 1384.
- 1278 Lozupone, C., Lladser, M.E., Knights, D., Stombaugh, J. and Knight, R. 2011. UniFrac: an
1279 effective distance metric for microbial community comparison. *ISME J* 5(2), 169-172.
- 1280 Luo, S., Nguyen, K.T., Nguyen, B.T.T., Feng, S., Shi, Y., Elsayed, A., Zhang, Y., Zhou, X., Wen,
1281 B., Chierchia, G., Talbot, H., Bourouina, T., Jiang, X. and Liu, A.Q. 2021. Deep learning-
1282 enabled imaging flow cytometry for high-speed Cryptosporidium and Giardia detection.
1283 *Cytometry Part A* 99(11), 1123-1133.

- 1284 Lytle, D.A., Pfaller, S., Muhlen, C., Struewing, I., Triantafyllidou, S., White, C., Hayes, S., King,
1285 D. and Lu, J. 2021. A comprehensive evaluation of monochloramine disinfection on water
1286 quality, *Legionella* and other important microorganisms in a hospital. *Water Research* 189,
1287 116656.
- 1288 Mahajna, A., Dinkla, I.J.T., Euverink, G.J.W., Keesman, K.J. and Jayawardhana, B. 2022. Clean
1289 and Safe Drinking Water Systems via Metagenomics Data and Artificial Intelligence: State-of-
1290 the-Art and Future Perspective. *Front Microbiol* 13, 832452.
- 1291 Mao, Y., Wang, X., Yang, H., Wang, H. and Xie, Y.F. 2014. Effects of ozonation on disinfection
1292 byproduct formation and speciation during subsequent chlorination. *Chemosphere* 117, 515-
1293 520.
- 1294 Mays, L.W. (2000) *Water Distribution System Handbook*, McGraw-Hill Education, New York.
- 1295 Mazhar, M.A., Khan, N.A., Ahmed, S., Khan, A.H., Hussain, A., Rahisuddin, Changani, F.,
1296 Yousefi, M., Ahmadi, S. and Vambol, V. 2020. Chlorination disinfection by-products in
1297 municipal drinking water – A review. *Journal of Cleaner Production* 273, 123159.
- 1298 Miake-Lye, I.M., Hempel, S., Shanman, R. and Shekelle, P.G. 2016. What is an evidence map?
1299 A systematic review of published evidence maps and their definitions, methods, and products.
1300 *Systematic reviews* 5, article No. 28.
- 1301 Moodley, T. and van der Haar, D. 2019 *Comparisons in Drinking Water Systems Using K-Means
1302 and A-Priori to Find Pathogenic Bacteria Genera*. Kim, K.J. and Baek, N. (eds), pp. 351-359,
1303 Springer Singapore, Singapore.
- 1304 Mostafavifar, M., Wertz, J. and Borchers, J. 2012. A systematic review of the effectiveness of
1305 kinesio taping for musculoskeletal injury. *The Physician and sportsmedicine* 40(4), 33-40.
- 1306 Mounce, S.R., Mounce, R.B. and Boxall, J.B. 2010. Novelty detection for time series data
1307 analysis in water distribution systems using support vector machines. *Journal of
1308 Hydroinformatics* 13(4), 672-686.
- 1309 Muharemi, F., Logofătu, D. and Leon, F. 2019. Machine learning approaches for anomaly
1310 detection of water quality on a real-world data set. *Journal of Information and
1311 Telecommunication* 3(3), 294-307.
- 1312 Mukhopadhyay, A., Duttagupta, S. and Mukherjee, A. 2022. Emerging organic contaminants in
1313 global community drinking water sources and supply: A review of occurrence, processes and
1314 remediation. *Journal of Environmental Chemical Engineering* 10(3), 107560.
- 1315 Narita, K., Matsui, Y., Matsushita, T. and Shirasaki, N. 2023. Screening priority pesticides for
1316 drinking water quality regulation and monitoring by machine learning: Analysis of factors
1317 affecting detectability. *Journal of Environmental Management* 326, 116738.
- 1318 Oh, S., Hossen, I., Luglio, J., Justin, G., Richie, J.E., Medeiros, H. and Lee, C.H. 2021. On-Site/In
1319 Situ Continuous Detecting ppb-Level Metal Ions in Drinking Water Using Block Loop-Gap
1320 Resonators and Machine Learning. *IEEE Transactions on Instrumentation and Measurement*
1321 70, 1-9.
- 1322 Olikier, N. and Ostfeld, A. 2014. Minimum volume ellipsoid classification model for
1323 contamination event detection in water distribution systems. *Environmental Modelling &
1324 Software* 57, 1-12.
- 1325 Ortiz-Lopez, C., Bouchard, C. and Rodriguez, M. 2022. Machine learning models with potential
1326 application to predict source water quality for treatment purposes: a critical review.
1327 *Environmental Technology Reviews* 11(1), 118-147.
- 1328 Pan, R., Zhang, T.-Y., Zheng, Z.-X., Ai, J., Ye, T., Zhao, H.-X., Hu, C.-Y., Tang, Y.-L., Fan, J.-J.,
1329 Geng, B. and Xu, B. 2023. Insight into mixed chlorine/chloramines conversion and associated

- 1330 water quality variability in drinking water distribution systems. *Science of The Total*
1331 *Environment* 880, 163297.
- 1332 Pandian, A.M.K., Rajamehala, M., Singh, M.V.P., Sarojini, G. and Rajamohan, N. 2022. Potential
1333 risks and approaches to reduce the toxicity of disinfection by-product—A review. *Science of the*
1334 *Total Environment* 822, 153323.
- 1335 Pang, G., Cao, L., Chen, L., Lian, D. and Liu, H. 2018. Sparse Modeling-Based Sequential
1336 Ensemble Learning for Effective Outlier Detection in High-Dimensional Numeric Data.
1337 *Proceedings of the AAAI Conference on Artificial Intelligence* 32(1).
- 1338 Park, J., Park, J.-H., Choi, J.-S., Joo, J.C., Park, K., Yoon, H.C., Park, C.Y., Lee, W.H. and Heo,
1339 T.-Y. 2020. Ensemble Model Development for the Prediction of a Disaster Index in Water
1340 Treatment Systems. *Water* 12(11), 3195.
- 1341 Park, M., Anumol, T. and Snyder, S.A. 2015. Modeling approaches to predict removal of trace
1342 organic compounds by ozone oxidation in potable reuse applications. *Environmental Science:*
1343 *Water Research & Technology* 1(5), 699-708.
- 1344 Peleato, N.M. 2022. Application of convolutional neural networks for prediction of disinfection
1345 by-products. *Scientific Reports* 12(1), 612.
- 1346 Peleato, N.M., Legge, R.L. and Andrews, R.C. 2018. Neural networks for dimensionality
1347 reduction of fluorescence spectra and prediction of drinking water disinfection by-products.
1348 *Water Research* 136, 84-94.
- 1349 Pifer, A.D. and Fairey, J.L. 2012. Improving on SUVA₂₅₄ using fluorescence-PARAFAC analysis
1350 and asymmetric flow-field flow fractionation for assessing disinfection byproduct formation
1351 and control. *Water Research* 46(9), 2927-2936.
- 1352 Pinto, A.J., Schroeder, J., Lunn, M., Sloan, W. and Raskin, L. 2014. Spatial-temporal survey and
1353 occupancy-abundance modeling to predict bacterial community dynamics in the drinking
1354 water microbiome. *mBio* 5(3), e01135-01114.
- 1355 Podgorski, J. and Berg, M. 2022. Global analysis and prediction of fluoride in groundwater. *Nat*
1356 *Commun* 13(1), 4232.
- 1357 Prescott, J.F. 2017. History and Current Use of Antimicrobial Drugs in Veterinary Medicine.
1358 *Microbiology Spectrum* 5(6), 10.1128/microbiolspec.arba-0002-2017.
- 1359 Proctor, C.R., Lee, J., Yu, D., Shah, A.D. and Whelton, A.J. 2020. Wildfire caused widespread
1360 drinking water distribution network contamination. *AWWA Water Science* 2(4), e1183.
- 1361 Qian, K., Jiang, J., Ding, Y. and Yang, S. 2020 Deep Learning Based Anomaly Detection in Water
1362 Distribution Systems, pp. 1-6.
- 1363 Ramos-Martínez, E., Herrera, M., Izquierdo, J. and Pérez-García, R. 2014. Ensemble of naïve
1364 Bayesian approaches for the study of biofilm development in drinking water distribution
1365 systems. *International Journal of Computer Mathematics* 91(1), 135-146.
- 1366 Ramos-Martínez, E., Herrera, M., Izquierdo, J. and Pérez-García, R. 2016. A Multi-disciplinary
1367 Procedure to Ascertain Biofilm Formation in Drinking Water Pipes. *International Congress on*
1368 *Environmental Modelling and Software* 13.
- 1369 Redondo-Hasselerharm, P.E., Cserbik, D., Flores, C., Farré, M.J., Sanchís, J., Alcolea, J.A., Planas,
1370 C., Caixach, J. and Villanueva, C.M. 2022. Insights to estimate exposure to regulated and
1371 non-regulated disinfection by-products in drinking water. *Journal of Exposure Science &*
1372 *Environmental Epidemiology*, 1-11.
- 1373 Renwick, D.V., Heinrich, A., Weisman, R., Arvanaghi, H. and Rotert, K. 2019. Potential Public
1374 Health Impacts of Deteriorating Distribution System Infrastructure. *J Am Water Works Assoc*
1375 111(2), 42-53.

- 1376 Richards, C.E., Tzachor, A., Avin, S. and Fenner, R. 2023. Rewards, risks and responsible
1377 deployment of artificial intelligence in water systems. *Nature Water* 1(5), 422-432.
- 1378 Richardson, S.D., Thruston, A.D., Caughran, T.V., Chen, P.H., Collette, T.W., Floyd, T.L., Schenck,
1379 K.M., Lykins, B.W., Sun, G.-r. and Majetich, G. 1999. Identification of New Ozone
1380 Disinfection Byproducts in Drinking Water. *Environmental Science & Technology* 33(19),
1381 3368-3377.
- 1382 Ritter, K.J., Carruthers, E., Carson, C.A., Ellender, R.D., Harwood, V.J., Kingsley, K., Nakatsu, C.,
1383 Sadowsky, M., Shear, B., West, B., Whitlock, J.E., Wiggins, B.A. and Wilbur, J.D. 2003.
1384 Assessment of statistical methods used in library-based approaches to microbial source
1385 tracking. *J Water Health* 1(4), 209-223.
- 1386 Roca, I., Akova, M., Baquero, F., Carlet, J., Cavaleri, M., Coenen, S., Cohen, J., Findlay, D.,
1387 Gyssens, I., Heuer, O.E., Kahlmeter, G., Kruse, H., Laxminarayan, R., Liébana, E., López-
1388 Cerero, L., MacGowan, A., Martins, M., Rodríguez-Baño, J., Rolain, J.M., Segovia, C.,
1389 Sigauque, B., Tacconelli, E., Wellington, E. and Vila, J. 2015. The global threat of
1390 antimicrobial resistance: science for intervention. *New Microbes New Infect* 6, 22-29.
- 1391 Rodriguez, M.J., Sérodes, J.B. and Levallois, P. 2004. Behavior of trihalomethanes and haloacetic
1392 acids in a drinking water distribution system. *Water Res* 38(20), 4367-4382.
- 1393 Rodriguez-Perez, J., Leigh, C., Liquet, B., Kermorvant, C., Peterson, E., Sous, D. and Mengersen,
1394 K. 2020. Detecting Technical Anomalies in High-Frequency Water-Quality Data Using
1395 Artificial Neural Networks. *Environmental Science & Technology* 54(21), 13719-13730.
- 1396 Sandler, M., Howard, A.G., Zhu, M., Zhmoginov, A. and Chen, L.-C. 2018 *MobileNetV2:*
1397 *Inverted Residuals and Linear Bottlenecks*, pp. 4510-4520.
- 1398 Sathasivan, A., Fisher, I. and Tam, T. 2008. Onset of severe nitrification in mildly nitrifying
1399 chloraminated bulk waters and its relation to biostability. *Water Research* 42(14), 3623-3632.
- 1400 Savadatti, M.B., Dhivya, M., Meghanashree, C., Navya, M., Lokesh, Y. and Kawri, N. 2022 An
1401 Overview of Predictive Analysis based on Machine learning Techniques, pp. 1-6.
- 1402 Senoro, D.B., de Jesus, K.L.M., Nolos, R.C., Lamac, M.R.L., Deseo, K.M. and Tabelin, C.B. 2022.
1403 In Situ Measurements of Domestic Water Quality and Health Risks by Elevated Concentration
1404 of Heavy Metals and Metalloids Using Monte Carlo and MLGI Methods. *Toxics* 10(7), 342.
- 1405 Shi, Y., Babatunde, A., Bockelmann-Evans, B., Li, Q. and Zhang, L. 2020. On-going nitrification
1406 in chloraminated drinking water distribution system (DWDS) is conditioned by hydraulics and
1407 disinfection strategies. *Journal of Environmental Sciences* 96, 151-162.
- 1408 Shi, Y., Wang, J., Wang, Q., Jia, Q., Yan, F., Luo, Z.-H. and Zhou, Y.-N. 2022. Supervised Machine
1409 Learning Algorithms for Predicting Rate Constants of Ozone Reaction with Micropollutants.
1410 *Industrial & Engineering Chemistry Research* 61(24), 8359-8367.
- 1411 Simpson, A.M.A. and Mitch, W.A. 2022. Chlorine and ozone disinfection and disinfection
1412 byproducts in postharvest food processing facilities: A review. *Critical Reviews in*
1413 *Environmental Science and Technology* 52(11), 1825-1867.
- 1414 Sincak, P., Ondo, J., Kaposztasova, D., Vircikova, M., Vranayova, Z. and Sabol, J. 2014. Artificial
1415 intelligence in public health prevention of legionellosis in drinking water systems. *Int J Environ*
1416 *Res Public Health* 11(8), 8597-8611.
- 1417 Singh, D., Vardhan, M., Sahu, R., Chatterjee, D., Chauhan, P. and Liu, S. 2023. Machine-learning-
1418 and deep-learning-based streamflow prediction in a hilly catchment for future scenarios using
1419 CMIP6 GCM data. *Hydrol. Earth Syst. Sci.* 27(5), 1047-1075.

- 1420 Singh, K.P. and Gupta, S. 2012. Artificial intelligence based modeling for predicting the
1421 disinfection by-products in water. *Chemometrics and Intelligent Laboratory Systems* 114, 122-
1422 131.
- 1423 Sluban, B. and Lavrač, N. 2015. Relating ensemble diversity and performance: A study in class
1424 noise detection. *Neurocomputing* 160, 120-131.
- 1425 Smith, A., Sterba-Boatwright, B. and Mott, J. 2010. Novel application of a statistical technique,
1426 Random Forests, in a bacterial source tracking study. *Water Res* 44(14), 4067-4076.
- 1427 Sokolova, M. and Lapalme, G. 2009. A systematic analysis of performance measures for
1428 classification tasks. *Information Processing & Management* 45(4), 427-437.
- 1429 Speight, V.L., Mounce, S.R. and Boxall, J.B. 2019. Identification of the causes of drinking water
1430 discolouration from machine learning analysis of historical datasets. *Environmental Science:
1431 Water Research & Technology* 5(4), 747-755.
- 1432 Srivastav, A.L., Patel, N. and Chaudhary, V.K. 2020. Disinfection by-products in drinking water:
1433 Occurrence, toxicity and abatement. *Environmental Pollution* 267, 115474.
- 1434 Steel, R.G.D. and Torrie, J.H. (1960) *Principles and procedures of statistics : with special reference
1435 to the biological sciences*, McGraw-Hill, New York.
- 1436 Sudhakaran, S. and Amy, G.L. 2013. QSAR models for oxidation of organic micropollutants in
1437 water based on ozone and hydroxyl radical rate constants and their chemical classification.
1438 *Water Res* 47(3), 1111-1122.
- 1439 Syafrudin, M., Kristanti, R.A., Yuniarto, A., Hadibarata, T., Rhee, J., Al-Onazi, W.A., Algarni, T.S.,
1440 Almarri, A.H. and Al-Mohaimed, A.M. 2021. Pesticides in drinking water—a review.
1441 *International Journal of Environmental Research and Public Health* 18(2), 468.
- 1442 Taheran, M., Naghdi, M., Brar, S.K., Verma, M. and Surampalli, R.Y. 2018. Emerging
1443 contaminants: here today, there tomorrow! *Environmental Nanotechnology, Monitoring &
1444 Management* 10, 122-126.
- 1445 Tian, C., Feng, C., Chen, L. and Wang, Q. 2020. Impact of water source mixture and population
1446 changes on the Al residue in megalopolitan drinking water. *Water Res* 186, 116335.
- 1447 Tolaymat, T.M., El Badawy, A.M., Genaidy, A., Scheckel, K.G., Luxton, T.P. and Suidan, M. 2010.
1448 An evidence-based environmental perspective of manufactured silver nanoparticle in
1449 syntheses and applications: a systematic review and critical appraisal of peer-reviewed
1450 scientific papers. *Science of the Total Environment* 408(5), 999-1006.
- 1451 U.S. EPA 2022 *Information about Public Water Systems*, EPA Office of Ground Water and
1452 Drinking Water, Washington, DC, U.S.A.
- 1453 Umesh C. Gupta, S.C.G. 2011. Heavy Metal Toxicity in Humans and its Preventive and Control
1454 Measures. *Current Nutrition & Food Science* 7(4), 221-231.
- 1455 Valko, M., Morris, H. and Cronin, M.T. 2005. Metals, toxicity and oxidative stress. *Curr Med
1456 Chem* 12(10), 1161-1208.
- 1457 Vinyals, O., Babuschkin, I., Czarnecki, W.M., Mathieu, M., Dudzik, A., Chung, J., Choi, D.H.,
1458 Powell, R., Ewalds, T., Georgiev, P., Oh, J., Horgan, D., Kroiss, M., Danihelka, I., Huang, A.,
1459 Sifre, L., Cai, T., Agapiou, J.P., Jaderberg, M., Vezhnevets, A.S., Leblond, R., Pohlen, T.,
1460 Dalibard, V., Budden, D., Sulsky, Y., Molloy, J., Paine, T.L., Gulcehre, C., Wang, Z., Pfaff, T.,
1461 Wu, Y., Ring, R., Yogatama, D., Wunsch, D., McKinney, K., Smith, O., Schaul, T., Lillicrap,
1462 T., Kavukcuoglu, K., Hassabis, D., Apps, C. and Silver, D. 2019. Grandmaster level in
1463 StarCraft II using multi-agent reinforcement learning. *Nature* 575(7782), 350-354.
- 1464 Wagner, E.D. and Plewa, M.J. 2017. CHO cell cytotoxicity and genotoxicity analyses of
1465 disinfection by-products: An updated review. *J Environ Sci (China)* 58, 64-76.

- 1466 Walesch, S., Birkelbach, J., Jézéquel, G., Haeckl, F.P.J., Hegemann, J.D., Hesterkamp, T., Hirsch,
1467 A.K.H., Hammann, P. and Müller, R. 2023. Fighting antibiotic resistance—strategies and
1468 (pre)clinical developments to find new antibacterials. *EMBO reports* 24(1), e56033.
- 1469 Wang, C., Yang, H., Liu, H., Zhang, X.-X. and Ma, L. 2023. Anthropogenic contributions to
1470 antibiotic resistance gene pollution in household drinking water revealed by machine-learning-
1471 based source-tracking. *Water Research* 246, 120682.
- 1472 Weigert, M., Schmidt, U., Haase, R., Sugawara, K. and Myers, G. 2020. Star-convex Polyhedra
1473 for 3D Object Detection and Segmentation in Microscopy, pp. 3655-3662, IEEE Computer
1474 Society.
- 1475 WHO (2011) Guidelines for drinking-water quality, World Health Organization, Geneva.
- 1476 WHO 2021. Antimicrobial resistance.
- 1477 Willard, J., Jia, X., Xu, S., Steinbach, M. and Kumar, V. 2022. Integrating Scientific Knowledge
1478 with Machine Learning for Engineering and Environmental Systems. *ACM Comput. Surv.*
1479 55(4), Article 66.
- 1480 Willmott, C.J. and Matsuura, K. 2005. Advantages of the mean absolute error (MAE) over the
1481 root mean square error (RMSE) in assessing average model performance. *Climate Research*
1482 30(1), 79-82.
- 1483 Wołos, A., Koszelewski, D., Roszak, R., Szymkuć, S., Moskal, M., Ostaszewski, R., Herrera, B.T.,
1484 Maier, J.M., Brezicki, G., Samuel, J., Lummiss, J.A.M., McQuade, D.T., Rogers, L. and
1485 Grzybowski, B.A. 2022. Computer-designed repurposing of chemical wastes into drugs.
1486 *Nature* 604(7907), 668-676.
- 1487 Wołowiec, M., Komorowska-Kaufman, M., Pruss, A., Rzepa, G. and Bajda, T. 2019. Removal of
1488 Heavy Metals and Metalloids from Water Using Drinking Water Treatment Residuals as
1489 Adsorbents: A Review. *Minerals* 9(8), 487.
- 1490 Wolpert, D.H. 1992. Stacked generalization. *Neural Networks* 5(2), 241-259.
- 1491 Wu, J., Song, C., Dubinsky, E.A. and Stewart, J.R. 2020. Tracking Major Sources of Water
1492 Contamination Using Machine Learning. *Front Microbiol* 11, 616692.
- 1493 Wu, Y., Jiang, P., Goh, S.G., Yu, K., Chen, Y., He, Y. and Gin, K.Y.H. 2022. Predicting Relative
1494 Risk of Antimicrobial Resistance using Machine Learning Methods. *IFAC-PapersOnLine*
1495 55(10), 1266-1271.
- 1496 Xiao, R., Ou, T., Ding, S., Fang, C., Xu, Z. and Chu, W. 2023. Disinfection by-products as
1497 environmental contaminants of emerging concern: a review on their occurrence, fate and
1498 removal in the urban water cycle. *Critical Reviews in Environmental Science and Technology*
1499 53(1), 19-46.
- 1500 Xie, Y., Sattari, K., Zhang, C. and Lin, J. 2023. Toward autonomous laboratories: Convergence
1501 of artificial intelligence and experimental automation. *Progress in Materials Science* 132,
1502 101043.
- 1503 Xu, X., Talbot, S. and Selvaraja, T. 2020. ParasNet: Fast parasites detection with neural networks.
1504 arXiv preprint arXiv:2002.11327.
- 1505 Zainurin, S.N., Wan Ismail, W.Z., Mahamud, S.N.I., Ismail, I., Jamaludin, J., Ariffin, K.N.Z. and
1506 Wan Ahmad Kamil, W.M. 2022. Advancements in monitoring water quality based on various
1507 sensing methods: a systematic review. *International Journal of Environmental Research and*
1508 *Public Health* 19(21), 14080.
- 1509 Zanoni, M.G., Majone, B. and Bellin, A. 2022. A catchment-scale model of river water quality
1510 by Machine Learning. *Science of the Total Environment* 838, 156377.

- 1511 Zhang, C., Brown, P.J.B. and Hu, Z. 2018. Thermodynamic properties of an emerging chemical
1512 disinfectant, peracetic acid. *Science of the Total Environment* 621, 948-959.
- 1513 Zhang, C., Hu, Z., Li, P. and Gajaraj, S. 2016. Governing factors affecting the impacts of silver
1514 nanoparticles on wastewater treatment. *Science of The Total Environment*.
- 1515 Zhang, C. and Lu, J. 2021a. Legionella: A Promising Supplementary Indicator of Microbial
1516 Drinking Water Quality in Municipal Engineered Water Systems. *Front Environ Sci* 9, 1-22.
- 1517 Zhang, C. and Lu, J. 2021b. *Legionella*: a supplementary indicator of microbial water quality in
1518 municipal engineered water systems. *Frontiers in Environmental Science* 9, 684319.
- 1519 Zhang, C. and Lu, J. 2021c. Optimizing disinfectant residual dosage in engineered water systems
1520 to minimize the overall health risks of opportunistic pathogens and disinfection by-products.
1521 *Science of The Total Environment* 770, 145356.
- 1522 Zhang, C., Struewing, I., Mistry, J.H., Wahman, D.G., Pressman, J. and Lu, J. 2021. Legionella
1523 and other opportunistic pathogens in full-scale chloraminated municipal drinking water
1524 distribution systems. *Water Res* 205, 117571.
- 1525 Zhang, Y., Gao, X., Smith, K., Inial, G., Liu, S., Conil, L.B. and Pan, B. 2019. Integrating water
1526 quality and operation into prediction of water production in drinking water treatment plants by
1527 genetic algorithm enhanced artificial neural network. *Water Research* 164, 114888.
- 1528 Zhong, S., Zhang, K., Bagheri, M., Burken, J.G., Gu, A., Li, B., Ma, X., Marrone, B.L., Ren, Z.J.,
1529 Schrier, J., Shi, W., Tan, H., Wang, T., Wang, X., Wong, B.M., Xiao, X., Yu, X., Zhu, J.-J. and
1530 Zhang, H. 2021. Machine Learning: New Ideas and Tools in Environmental Science and
1531 Engineering. *Environmental Science & Technology* 55(19), 12741-12754.
- 1532 Zhou, Q., Bian, Z., Yang, D. and Fu, L. 2023. Stability of Drinking Water Distribution Systems
1533 and Control of Disinfection By-Products. *Toxics* 11(7), 606.
- 1534 Zhou, X., Tang, Z., Xu, W., Meng, F., Chu, X., Xin, K. and Fu, G. 2019. Deep learning identifies
1535 accurate burst locations in water distribution networks. *Water Research* 166, 115058.
- 1536 Zhu, M., Wang, J., Yang, X., Zhang, Y., Zhang, L., Ren, H., Wu, B. and Ye, L. 2022. A review of
1537 the application of machine learning in water quality evaluation. *Eco-Environment & Health*
1538 1(2), 107-116.
- 1539 Ziegelbauer, K., Speich, B., Mäusezahl, D., Bos, R., Keiser, J. and Utzinger, J. 2012. Effect of
1540 sanitation on soil-transmitted helminth infection: systematic review and meta-analysis. *PLOS*
1541 *Medicine* 9(1), article No. e1001162.
- 1542

Using Machine Learning to Ensure Drinking Water Quality

Water Source

- Source tracking
- Water source quality
- ...

Water Utility

- Treatment efficiency
- Disinfectant dosage
- ...

Distribution System

- Opportunistic pathogens
- Heavy metals
- ...

Premise Plumbing

- Disinfection byproducts
- Opportunistic pathogens
- ...

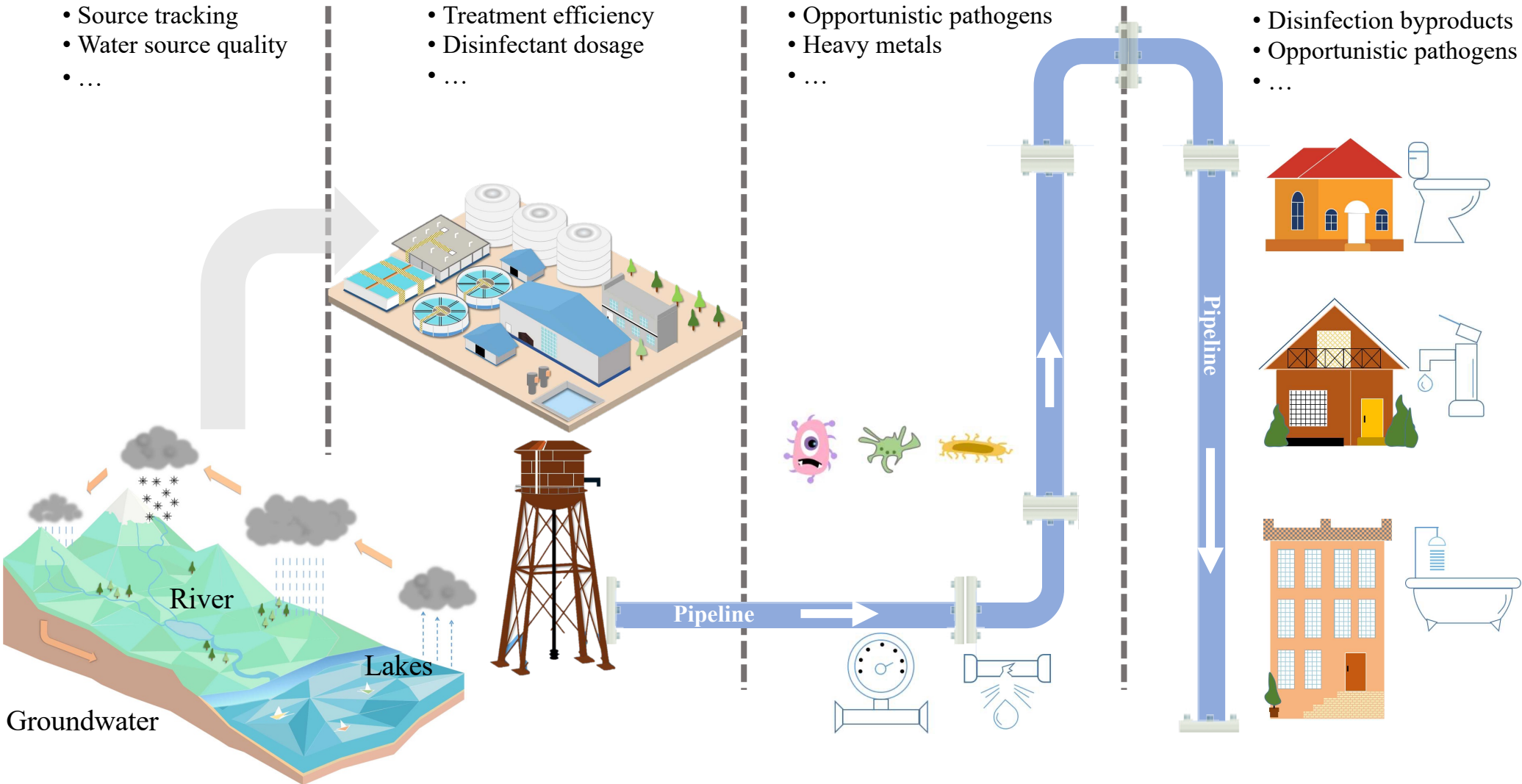


Figure 1. A stage-by-stage overview of the application of machine learning to ensure drinking water quality

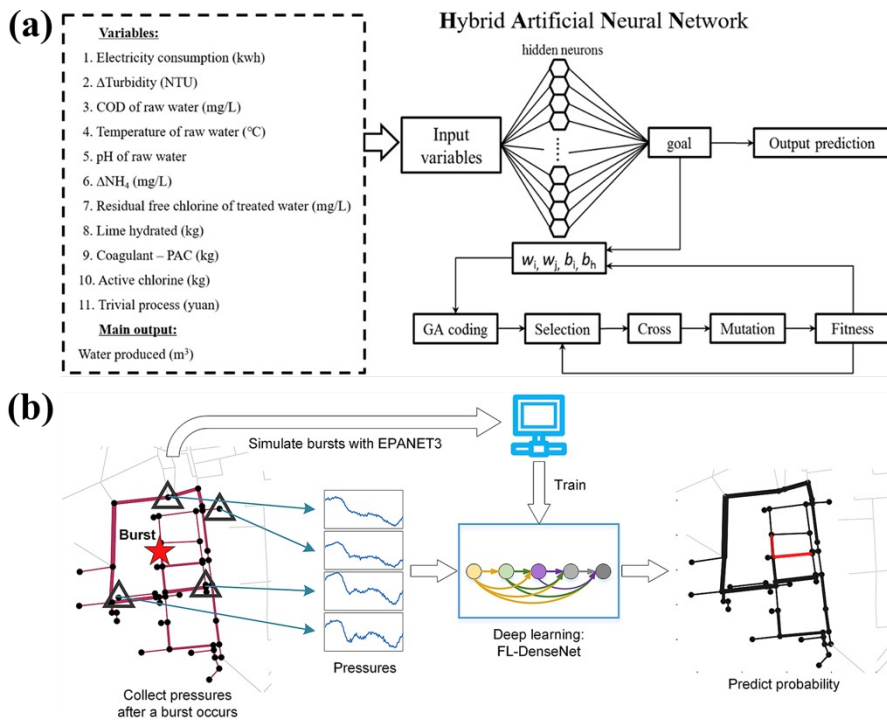


Figure 2. (a) Input and output variables used for modeling and the proposed hybrid artificial neural network framework. Reproduced with permission from Zhang et al., 2019. Copyright 2019 Elsevier. (b) Schematic of fully-linear DenseNet (BLIFF) model for accurate identification of burst locations in EWS networks. Reproduced with permission from Zhou et al., 2019. Copyright 2019 Elsevier.

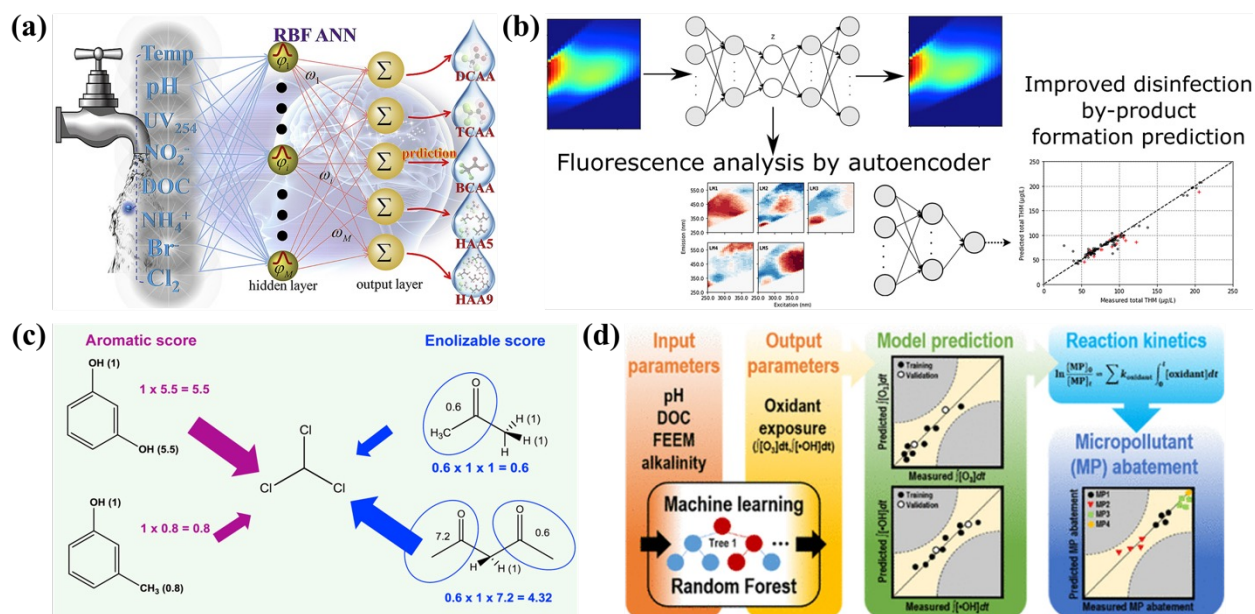


Figure 3. (a) Schematic of radial basis function (RBF) artificial neural network (ANN) model for prediction of disinfection by-products (DBPs). Reproduced with permission from Lin et al., 2020, Copyright 2020 Elsevier. (b) Schematic of parallel factor analysis (PARAFAC) model for prediction of DBPs. Reproduced with permission from Peleato et al., 2018, Copyright 2018 Elsevier. (c) Schematic of multiple linear regression (MLR) model for prediction of DBPs from organic precursors. Reproduced with permission from Bond and Graham 2017, Copyright 2017 Elsevier. (d) Schematic of random forest (RF) model for prediction of micropollutant abatement. Reproduced with permission from Cha et al., 2021, Copyright 2021 ACS.

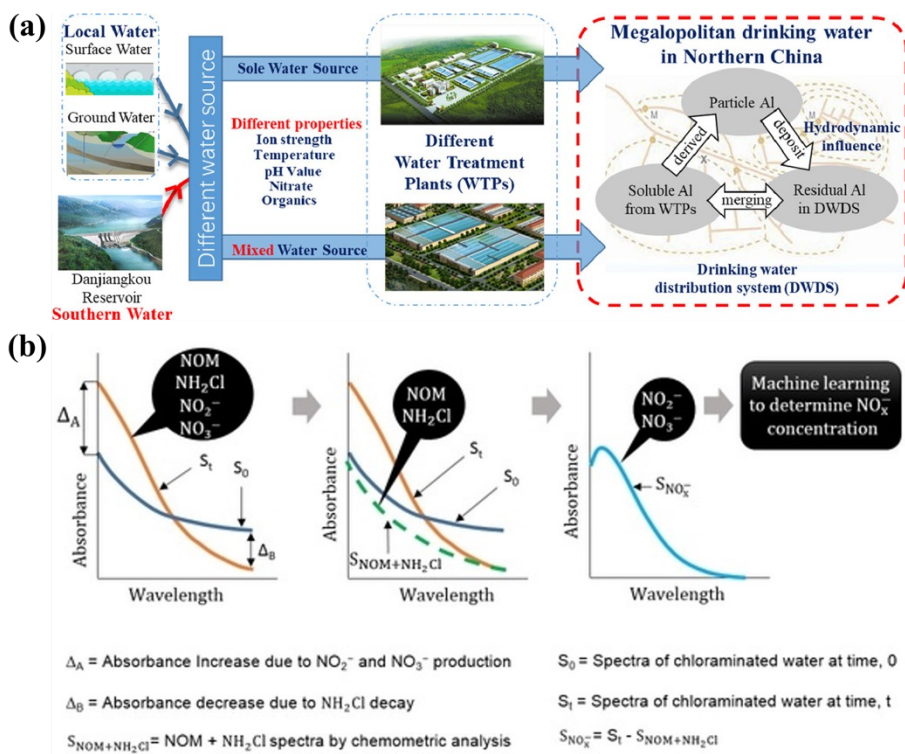


Figure 4. (a) The impact of water source mixture and population changes on the aluminum (Al) residue.

Reproduced with permission from Tian et al., 2020, Copyright 2020 Elsevier. (b) Prediction of nitrate and

nitrite concentrations over support vector regression (SVR) model. Reproduced with permission from

Hossain et al., 2021, Copyright 2021 MDPI.

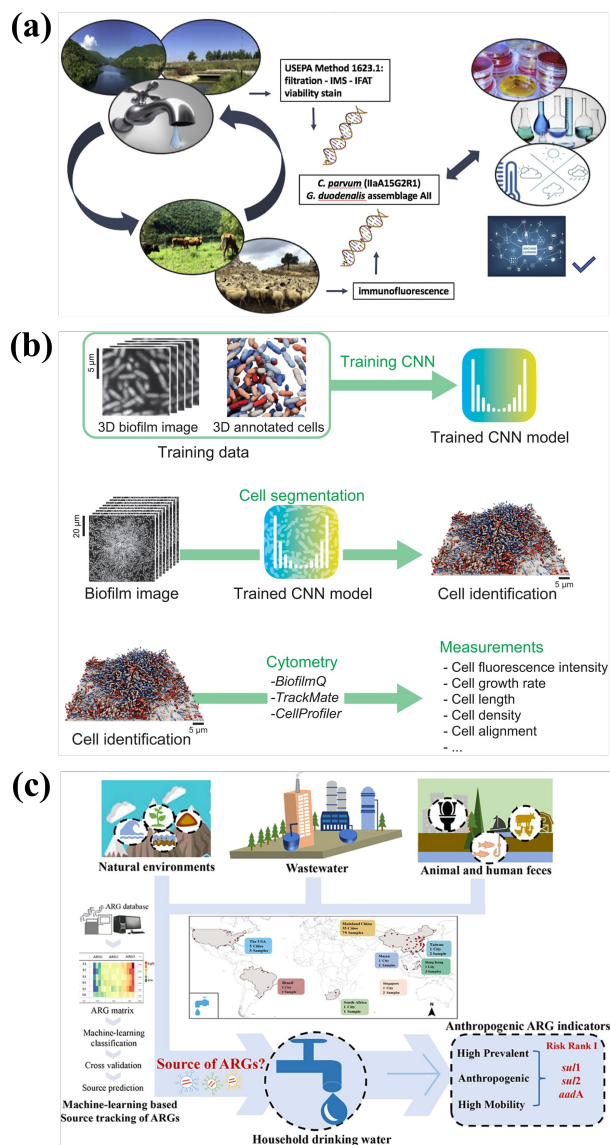


Figure 5. (a) The application of machine learning to predict contamination of *Cryptosporidium* and *Giardia* in surface water and drinking water. Reproduced with permission from Ligda et al., 2020, Copyright 2020 Elsevier. (b) Deep-learning-based workflow for single-cell measurements in three-dimensional biofilms. Reproduced with permission from Jelli et al., 2023, Copyright 2023 Elsevier. (c) SourceTracker was performed to investigate the pollution sources of antimicrobial resistance genes (ARGs) in household drinking water. Reproduced with permission from Wang et al., 2023, Copyright 2023 Elsevier.

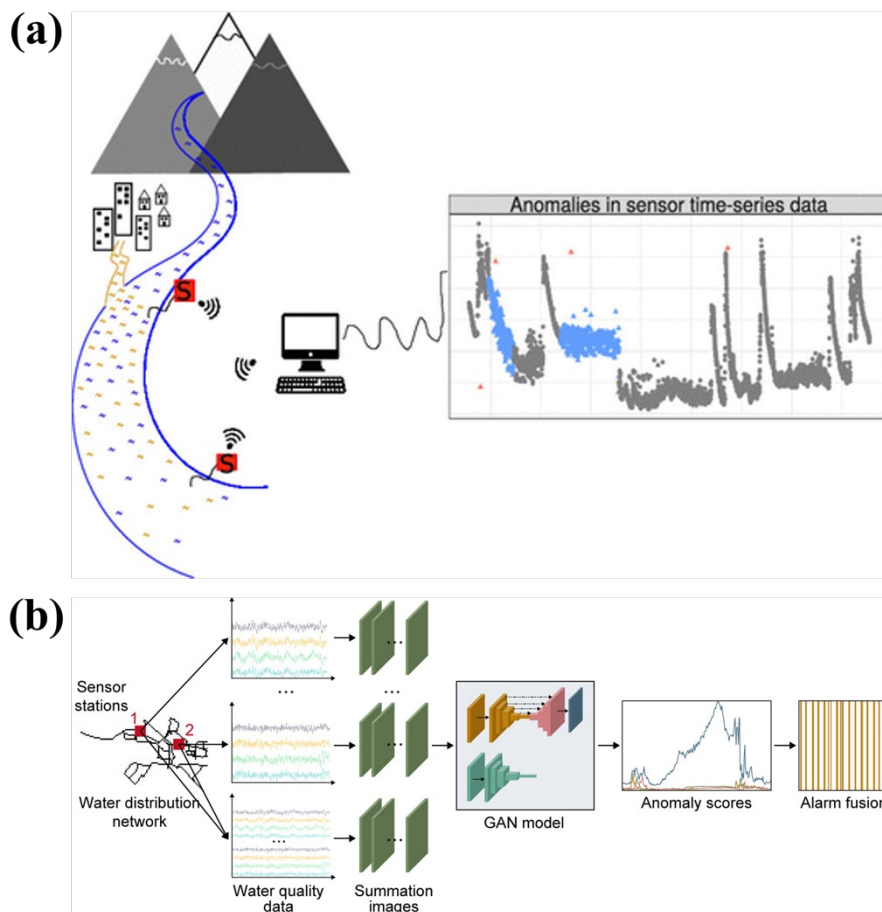


Figure 6. (a) Detection of technical anomalies in water quality using artificial neural network (ANN) model. Reproduced with permission from Rodriguez-Perez et al., 2020, Copyright 2020 ACS. (b) Detection of contamination events using generative adversarial network (GAN) model. Reproduced with permission from Li et al., 2023, Copyright 2023 Elsevier.

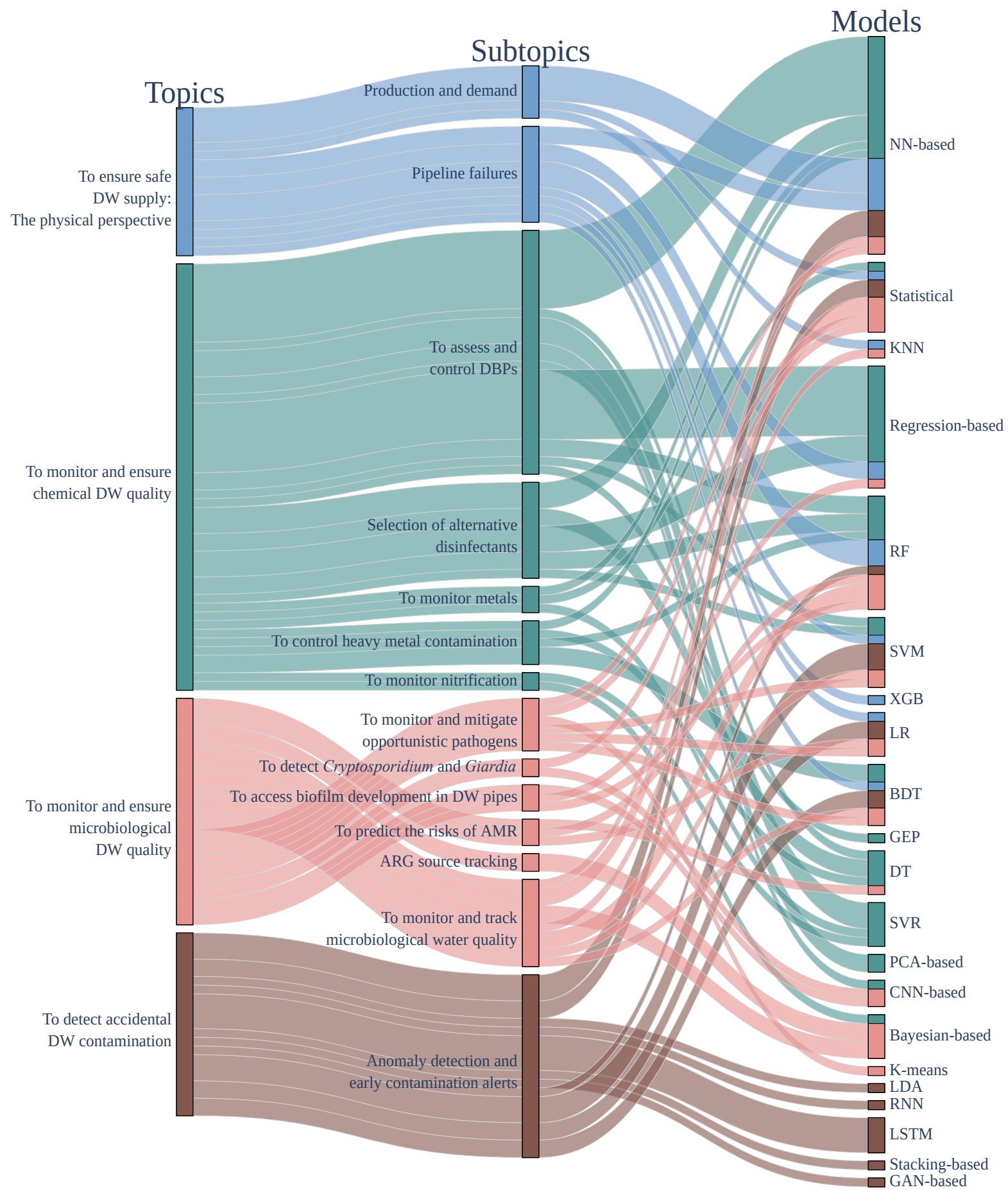


Figure 7. Distribution of machine learning models across drinking water distribution system research topics

Abbreviations: NN, Neural network; KNN, K-nearest neighbor; RF, random forest; SVM, support vector machine; XGB, extreme gradient boosting; LR, logistic regression; BDT, boosting decision tree; GEP, gene expression programming; DT, decision tree; SVR, support vector regression; PCA, principal component analysis; CNN, convolutional neural network; LDA, linear discriminant analysis; RNN, recurrent neural network; LSTM, long short-term memory; GAN, generative adversarial network; DW, drinking water.

Table 1 The applications of machine learning to ensure safe drinking water supply from the physical perspective

Topic	Task	Model	Inputs	Outputs	Metrics (Selected)	Performances (Selected)	Reference
Drinking water production and demand	Water production prediction	GA-ANN and ML-ANN	T, COD, and operational parameters	Water production of DWTPs	MSE, R^2 , MAPE	GA-ANN $R^2=0.93$ > ML-ANN	Zhang et al., 2019
	Short-term water demand prediction	GRUN, ANN, and SARIMA	Historical water demand data	15-min and 24-h prediction of water demand	MAE, MAPE, RMSE, NSE	GRUN > ANN and SARIMA	Guo et al., 2018
	Water demand forecasting	DAN2, FTDNN, and KNN	Daily water production and monthly water consumption	Daily, weekly, and monthly water demands	MAPE, accuracy, R^2 , MSE, and SSE	DAN2 accuracies: 96% to 98%	Ghiassi et al., 2017
Drinking water pipeline failures	Pipe burst localization	FL-DenseNet	Pressure measurements	Burst occurring likelihood per pipe	Accuracy	62.35% to 98.58%	Zhou et al., 2019
	Pipe failure prediction	AdvaML, Cox-pH, SRF, and SSVM	Pipe data and climate data	Failure/Hazard Index	C-index	AdvaML ≥ 0.8 > Cox-pH, SRF, and SSVM	Almheiri et al., 2021
	Disaster index prediction on WTS	RF and XGB	Facility specification and operational data	Disaster index	RMSE and R^2	XGB $R^2 = 0.86$ > RF	Park et al., 2020
	Earthquake damage prediction	RR, LR, BRT, and RF	Earthquake-related variables and pipe attributes	Binary classification of damage status	TE, TEP, RMSE, MAE, MASE, MPSE, SN, SP, TSS, and AUC	BRT > RR, LR, and RF in overall performance	Bagriacik et al., 2018

GA-ANN, Artificial neural network with genetic algorithm; ML-ANN, multi-layered artificial neural network; T, temperature; COD, chemical oxygen demand; DWTPs, drinking water treatment plants; MSE, mean squared error; R^2 , coefficient of determination; MAPE, mean absolute percentage error; GRUN, gated recurrent unit network; SARIMA, seasonal autoregressive integrated moving average; MAE, mean absolute error; RMSE, root-mean square error; NSE, Nash-Sutcliffe model efficiency; DAN2, dynamic artificial neural network; FTDNN, focused time-delay neural network; KNN, K-nearest neighbor; SSE, summing the squared differences; FL-DenseNet, fully-linear DenseNet; AdvaML, advanced meta-learning; Cox-pH, cox-proportional hazards; SRF, random survival forest; SSVM, survival support vector machine; C-index, concordance index; WTS, water treatment system; RF, random forest; XGB, extreme gradient boosting (XGBoost); RR, repair rate; LR, logistic regression; BRT, boosted regression trees; TE, error in total count; TEP, percentage error in total count; MASE, median absolute suburb error; MPSE, Median percentage suburb error; SN, sensitivity; SP, specificity; TSS, true skill statistics; AUC, area under the receiver operating characteristic (ROC) curve.

Table 2 The applications of machine learning to monitor and ensure chemical drinking water quality

Topic	Task	Model	Inputs	Outputs	Metrics (Selected)	Performances (Selected)	Reference
To assess and control DBPs	DBPs formation prediction	ANN, SVM, and GEP	pH, T, C_{Br^-} , $C_{Cl_2/DOC}$, t	C_{THMs}	MSE, RMSE and R^2	SVM > ANN and GEP	Singh and Gupta 2012
		Linear/log linear, and RBF-ANN	pH, T, UV_{254} , C_{DOC} , C_{Br^-} , $C_{residual_cl}$, $C_{NO_2^- - N}$, and $C_{NH_4^+ - N}$	C_{HAAs}	Accuracy, AAE	RBF-ANN > linear/log linear	Lin et al., 2020
		Linear/log linear, and RBF-ANN	pH, T, UV_{254} , C_{DOC} , C_{Br^-} , $C_{residual_free_cl}$, $C_{NO_2^- - N}$, and $C_{NH_4^+ - N}$	C_{THMs}	Accuracy and r_p	RBF ANN > linear/log linear	Hong et al., 2020
		Linear/log linear BP-ANN, and RBF-ANN		C_{HKs}	R^2	RBF ANN:0.799 > BP ANN and linear/log linear	Deng et al., 2021
		DTB	C_{NH_2Cl} , $C_{NHCl_2 + OC}$, pH, TDN, $C_{NO_2^- - N}$, TOC, and $C_{NH_4^+ - N}$	C_{THM4} and C_{HAAs}	R^2 and MSE	C_{THM4} : $R^2 = 0.56$ C_{HAAs} : $R^2 = 0.65$	Pan et al., 2023
		MLR, NN, RF, GPR and SVR	T, $C_{residual_cl}$, DOC, Turb, pH, Leit, and UV_{254}	C_{THMs} , C_{HAAs} , C_{DCAN} , C_{CPK} , and C_{TCP}	MSE	SVR, GPR > NN > RF > MLR	Hu et al., 2023
	Spectroscopic detection of DBPs	AE-NN, AE, PCA, and PARAFAC	Fluorescence spectra	C_{THMs} and C_{HAAs}	MAE, MSE	AE-NN > AE > PCA > PARAFAC	Peleato et al., 2018
		MLP, CNN, PARAFAC-MLP, PCA-MLP, and 3-way PLS		C_{THMs} , C_{HAAs} , and C_{TCMs}		CNN > MLP, PARAFAC-MLP, PCA-MLP, and 3-way PLS	Peleato 2022
		SVR-(linear, polynomial, RBF, or sigmoid)	UV-Vis spectra	C_{NH_2Cl}	R^2 and RMSE	SVR-RBF > SVR-Polynomial > SVR-linear and SVR-Sigmoid	Hossien et al., 2020
	DBPs formation mechanism analysis	MLR	Chemical descriptors	THM yield	R^2 and RMSE	$R^2 = 0.91$	Bond and Graham 2017
RF, SVR-RBF, SVR-linear, MLP, and MLR		Chemical descriptors	HAAs formation potential		RF > SVR-RBF, SVR-linear, MLP and MLR	Cordero et al., 2021	

DBPs, Disinfection by-products; GEP, gene expression programming; C_{Br^-} , Br concentration; $C_{Cl_2/DOC}$, dissolved organic carbon normalized chlorine dose; t , contact time; C_{THMs} , trihalomethane concentration; linear/log Linear, linear/log linear regression models; RBF-ANN, radial basis function ANN; UV_{254} , ultraviolet absorbance at 254 nm; C_{DOC} , dissolved organic carbon concentration; $C_{residual_cl}$, residue chlorine concentration; $C_{NO_2^- - N}$, nitrite concentration; $C_{NH_4^+ - N}$, ammonia concentration; C_{HAAs} , haloacetic acids concentration; AAE, average absolute error; $C_{residual_free_cl}$, residual free chlorine concentration; r_p , regression coefficients; BP-ANN, back propagation ANN; C_{HKs} , haloketones concentration; DTB: decision tree boost; C_{NH_2Cl} , monochloramine concentration; $C_{NHCl_2 + OC}$, dichloramine and organic chloramines concentration; TDN, total dissolved nitrogen; TOC, total organic carbon;

Table 2 The applications of machine learning to monitor and ensure chemical drinking water quality

MLR, multiple linear regression; GPR, Gaussian process regression; SVR, support vector regression; Turb, turbidity; Leit, electric conductivity of the water; C_{DCAN} , dichloroacetonitrile concentration; C_{CPK} , chloropicrin concentration; C_{TCP} , trichloropropanone concentration; AE-NN, autoencoder-neural network; PCA, principal component analysis; PARAFAC, parallel factors analysis; MLP, multi-layer perceptron network; CNN, convolutional neural network; 3-way PLS, 3-way partial least squares; C_{TCMS} , trichloromethane concentration; C_{NH_2Cl} , monochloramine concentration; MLP, multilayer perceptron.

Table 2 The applications of machine learning to monitor and ensure chemical drinking water quality (cont.)

Topic	Task	Model	Inputs	Outputs	Metrics (Selected)	Performances (Selected)	Reference
To select alternative disinfectants	Prediction of bromate formation by ozonation	MLR and ANN	C_{τ} , pH, $C_{BrO_3^-}$, T, UV, DOC, Alk, and $C_{NH_4^+-N}$	$C_{BrO_3^-}$	R^2	ANN = 0.98 > MLR	Legube et al., 2004
	Prediction of MP/organic contaminant abatement during ozonation	RF	pH, Alk, DOC, and FEEM	Oxidant exposures	R^2 and RMSE	$R^2 : 0.904$	Cha et al., 2021
		MLR, SVM, DT, RF, and DNN	NI _a , E_{LUMO} , E_{HOMO} , and $C_{Ene_val,min}$	$\log k_{O_3}$	R^2 , MSE, MAE and Q_{ext^2}	RF: $R^2 = 0.9113$	Shi et al., 2022
		DTB and SDT	k_{O_3} model: AMR, minHBa, n_X , and MDEC-24 ; k_{SO_4} model: AMR, SssO, and meanI	k_{O_3} and k_{SO_4}	R^2 and RMSE	DTB: $R^2 > 0.97$	Gupta and Basant 2016
	Estimation of the TOrcs removal	MLR, ANN, and PC-ANN	C_{O_3} , TOC, $k_{O_3,TOrc}$ and $k_{OH,TOrc}$	TOrcs removal	R^2 and RMSE	PC-ANN: $R^2 = 0.934$	Park et al., 2015
To monitor HM	Pb ions concentration detection	SVR	S_{11}	Pb concentration	RMS	0.71	Oh et al., 2021
	Spatial HM concentration mapping	MLGI (NN - PSO + EBK)	Geographical coordinates	Spatial concentration maps	MSE and r	$r \approx 1.0$	De Jesus et al., 2021
	Temporal-spatial map generating	Kriging interpolation	Spatial and temporal data	Temporal-spatial distribution of residual Al	-	-	Tian et al., 2020
To control HM contamination	As adsorption removal prediction	LightGBM, XGB, GBDT, and RF	adsorbent dosage, t , C_{As_init} , pH, T, A_{MOFs} , and N_{anions}	Adsorptive removal of As(V)	AAPRE, RMSE and R^2	LightGBM > XGBoost > GBDT > RF	Abdi et al., 2022
	HM removal prediction	MLP-ANN and RBF-ANN	adsorbent dosage, τ , and pH_{init}	Al, Cd, Co, Cu, Fe, and Pb ions removal efficiency	MSE and R^2	RBF-ANN > MLP-ANN	Hamidian et al., 2019
To monitor nitrification	Nitrification episodes classification	NB	16S rRNA profiling	Nitrification episodes: stable or failure	AUC	0.83	Gomez-Alvarez et al., 2020
	Estimate NOx concentrations	SVR	NOx absorbances at various wavelengths	C_{NO_x}	RMSE and R^2	RMSE < 0.04	Hossain et al., 2021

C_{τ} , Disinfectant concentration and contact time product; Alk, alkalinity; $C_{BrO_3^-}$, bromate concentration; FEEM, fluorescence excitation–emission matrix; DT, decision tree; DNN, deep neural network; NI_a, norm descriptors; E_{LUMO} and E_{HOMO} , energy of the lowest unoccupied molecular orbital and energy of the highest occupied molecular orbital; $C_{Ene_val,min}$, minimum valence shell orbital energy on carbon atom; Q_{ext^2} , external validation parameter; SDT, single decision tree; k_{O_3} and k_{SO_4} , the rate constants for the reactions of O_3 and SO_4^- respectively; AMR, antimicrobial resistance; minHBa, minimum E-states for (strong) hydrogen bond acceptors; n_X , number of halogen atoms; MDEC-24, molecular distance edge between all secondary and quaternary carbons, SSSC, sum of atom-type E-state; O^- , meanI, mean intrinsic state values.

Table 2 The applications of machine learning to monitor and ensure chemical drinking water quality (*cont.*)

TOrCs, trace organic compounds; PC-ANN, principal component ANN; C_{O_3} , applied ozone dose; $k_{O_3,TOrc}$ and $k_{OH,TOrc}$, rate constants of O_3 and $\cdot OH$ of TOrCs; HM, heavy metal; S_{11} , reflection coefficient; MLGI (NN-PSO+EBK), machine learning and geostatistical interpolation (neural network with the particle swarm optimization and empirical Bayesian kriging); r , Pearson's correlation coefficient; LightGBM, light gradient-boosting machine; GBDT, gradient boosting decision tree; t , contact time; C_{As_init} , initial arsenic concentration; A_{MOFs} , metal-organic frameworks surface area; N_{anions} , presence of anions; AAPRE, average absolute percent relative error; pH_{init} , initial pH; NB, naïve Bayes; AUC, area under the curve; NOx, nitrite and nitrate.

Table 3 The applications of machine learning to monitor and ensure microbiological drinking water quality

Topic	Task	Model	Inputs	Outputs	Metrics (Selected)	Performances (Selected)	Reference
To monitor and mitigate opportunistic pathogens	To simulate conditions for preventing legionleosis outbreak	NARA	Q and T	T profile of the water tank	Accuracy	>97%	Sincak et al., 2014
	Bacterium clustering	K-means	16S rRNA profiling	Clusters of bacteria	-	-	Moodley and Haar 2019
	Spatio-temporal clustering of high-risk; serogroup and contamination levels prediction	SaTScan, XGB, LR, and SVM	Survey, spatial and meteorological info., and risk level to <i>Legionella</i> ;	High-risk level clusters; serogroup of a sample and the contamination level	Accuracy and F1-score	XGBoost > SVM > LR	Brunello et al., 2022
To detect <i>Cryptosporidium</i> and <i>Giardia</i>	Image classification of <i>Cryptosporidium</i> and <i>Giardia</i> morphology	CNN	Cell level scattering image	Classification of <i>Cryptosporidium</i> , <i>Giardia</i> , or others	Accuracy	Accuracy: 95.6% for <i>Cryptosporidium</i> and 99.5% for <i>Giardia</i>	Xu et al., 2020
	<i>Cryptosporidium</i> and <i>Giardia</i> contamination intensity prediction	LDFA	Microbiological, physicochemical, and meteorological parameters	Multiple classification or binary classification	Accuracy, precision, recall, and F1-score	Accuracy > 99.6%	Luo et al., 2021
				(oo)cyst concentrations of <i>Cryptosporidium</i> and <i>Giardia</i>	Accuracy	Accuracy: 75% for <i>Cryptosporidium</i> and 69% for <i>Giardia</i>	Ligda et al., 2020
to access biofilm development in drinking water pipes	Biofilm development analysis	RT and RF	System physical and hydraulic characteristics, sampling and incubation, and physico-chemical of water	HPC	R	RF: 0.898	Ramos-Martínez et al., 2016
	Single-cell segmentation in 3D biofilms	StarDist OPP (CNN-based)	3D biofilm image	Cell identification	Precision and OSA	OSA = 3% Precision depends on IoU threshold	Jelli et al., 2023

NARA, Neural network designed on approximate reasoning architecture; Q , flow rate; LDFA, linear discriminant function analysis; RT, regression trees; HPC, heterotrophic plate count; OSA, over-segmentation abundances; IoU, intersection-over-union.

Table 3 The applications of machine learning to monitor and ensure microbiological drinking water quality (cont.)

Topic	Task	Model	Inputs	Outputs	Metrics (Selected)	Performances (Selected)	Reference
To predict risk of AMR and track sources of ARGs	Relative risk of AMR prediction	LR, DT, and RF	T, pH, ORP, EC, ρ , TDS, Sal, P, DO, Turb, and 24h rainfall	Relative risk score	Accuracy, precision, recall, F1-score and AUC	RF: AUC = 0.88 > DT, LR	Wu et al., 2022
	ARG source tracking	Bayesian-based	Metagenomic signatures of ARGs and microbial taxa	Relative contributions of ARGs	-	-	Chen et al., 2019
			Broad-spectrum ARG profiles	Proportion of pollution sources of AGGs	r	$r = 0.98$	Wang et al., 2023
To monitor and track microbiological water quality	Water source tracking	Bayesian-based	rep-PCR and ARA	Source membership	-	RMSEp < RMSEc	Ritter et al., 2003
			ARA	Source distribution	RMSE		Greenberg et al., 2010
			Bacterial 16S ribosomal RNA gene sequences	Source proportion	R^2	≥ 0.8	Knights et al., 2011
	Microbial contamination prediction	RF	ARA	Source classification	ARCC	82.3%	Smith et al., 2010
		XGB, KNN, NB, SVM, NN and RF	Weather, hydrologic and land cover data	Source classification	Accuracy and AUC	XGBoost > RF > KNN > NN > SVM > NB	Wu et al., 2020
		Hidden features of bacterial communities unveiling	Alpha and Beta diversity analyses	Sequencing data of the bacterial community	Clustering properties of bacterial community	Unweighted UniFrac score	-
UniFrac	-		-	Unweighted/weighted UniFrac score	-	Lozupone et al., 2011; Bruno et al., 2018; Ling et al., 2018; Li et al., 2017	

ORP, oxidation-reduction potential; EC, electrical conductance; ρ , resistivity; TDS, total dissolved solids; Sal, salinity; P, pressure; DO, dissolved oxygen; 24h rainfall, 24h accumulated rainfall; rep-PCR, repetitive element polymerase chain reaction; ARA, antibiotic resistance analysis; RMSEp, RMSE for posterior probability averaging estimator; RMSEc, RMSE for classification method estimator; ARCC, average rates of correct classification; NN, neural network.

Table 4 The applications of machine learning to detect accidental drinking water contamination

Topic	Task	Model	Inputs	Outputs	Metrics (Selected)	Performances (Selected)	Reference
To detect anomalies and contamination events in DW	Anomaly event detection	LR, LDA, SVM, ANN, DNN, RNN, and LSTM	T, C_{ClO_2} , pH, Redox, Leit, Turb, and Q	Event (Boolean)	F1-score	SVM: F1-score = 0.36	Muharemi et al., 2019
		LSTM				LSTM: F1-score = 0.80	Fehst et al., 2018
		LR, RF, XGB, xgbDART, and LSTM				LSTM: F1-score = 0.78	Qian et al., 2020
		LSTM and ARIMA	Turb and Leit		b.Acc, F1-score, and MCC	LSTM > ARIMA	Rodriguez-Perez et al., 2020
		Stacking-based and ANN	Cl ₂ , pH, Leit, T, TOC, and Turb		F1-score, R ² , and MSE	Stacking > ANN	Li et al., 2022
	GAN-based and MVE-based			FAR, F1-score, and EDR	GAN > MVE	Li et al., 2023	
	Contamination event detection	SVM	Cl ₂ , EC, pH, T, TOC, and Turb	Three-class-event classification	Accuracy and EDR	Accuracy: 0.83-0.97	Oliker and Ostfeld, 2014
	DW classification: potable vs. contaminated	SVM	UV-absorbance readings	Contamination event	Confusion matrix	False alarm: 0.19	Asheri Arnon et al., 2019

LDA, linear discriminant analysis; RNN, recurrent neural network; LSTM, long short-term memory; C_{ClO_2} , chlorine dioxide concentration; Redox, redox potential; xgbDART, extreme gradient boosting with dropouts meet multiple additive regression trees; ARIMA, auto-regressive integrated moving average; b.Acc, balanced accuracy; MCC, Matthews correlation coefficient; Cl₂, total chlorine; GANs, generative adversarial networks; MVE, minimum volume ellipsoid; FAR, false alarm rate; EDR, event detection rate.