

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

2 Feng Feng^{1,2†}, Yuanxun Zhang^{1†}, Zhenru Chen³, Jianyuan Ni⁴, Yuan Feng⁵, Yunchao Xie^{6*},

3 Chiqian Zhang^{7*}

4

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

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

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

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

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

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

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

18

19 †These authors contributed equally

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

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

41 **Keywords:** Drinking water quality; Engineered water systems; Artificial intelligence;
42 Opportunistic pathogens; Disinfection byproducts; Heavy metals

43

44 1. Introduction

45 Clean and safe drinking water is vital to public health and socioeconomic development ADDIN
46 EN.CITE (. Public water systems are essential drinking water sources in modern society. For
47 instance, in the United States (U.S.), over 90% of people obtain drinking water from approximately
48 150,000 public water systems (U.S. EPA, 2023). Public water systems provide water for human
49 consumption through engineered water systems, including drinking water distribution systems
50 (DWDSs) and building premise plumbing systems (Zhang and Lu, 2021b). Poor drinking water
51 quality causes disease outbreaks and chronic diseases, leading to significant socioeconomic losses
52 (Benedict et al., 2017; Craun et al., 2010; Lee et al., 2023). Therefore, assessing and ensuring water
53 quality, especially in engineered water systems, is critical to public health and the welfare of
54 society (WHO, 2011).

55
56 Assessing and ensuring water quality in engineered water systems and related settings are complex
57 (Li and Wu, 2019). Multiple ever-changing variables affect drinking water quality such as the
58 quality of water sources, treatment processes and technologies, water pipe materials, distribution
59 system configuration and length, natural disasters, and water stagnation in the pipes (Delpla et al.,
60 2009; Li and Wu, 2019; Proctor et al., 2020). For instance, drinking water effluent at water utilities
61 may be high but could deteriorate in engineered water systems because of microbial (re)growth
62 and inactivation, the formation of disinfection by-products (DBPs), pipe failures (e.g., breaks and
63 leaks), and the detachment of heavy metals from pipes (Liu et al., 2016). Meanwhile, drinking
64 water in engineered water systems and other settings have various hazardous agents such as
65 harmful microbes (especially opportunistic pathogens or OPs) (Li et al., 2019; Zhang and Lu,
66 2021a), DBPs (Benítez et al., 2021; Lee et al., 2013), heavy metals (Chowdhury et al., 2016;

67 Gonzalez et al., 2013), pesticides and herbicides (Syafudin et al., 2021) (Mukhopadhyay et al.,
68 2022), and other emerging contaminants (e.g., antimicrobials and microplastics) (Gogoi et al.,
69 2018; Kirstein et al., 2021; Taheran et al., 2018). Those agents are interconnected, and controlling
70 only one group of the agents frequently fails to secure drinking water quality. For instance,
71 increasing disinfectant residual concentrations in engineered systems suppresses microbial
72 (re)growth but promotes the formation of DBPs (Zhang and Lu, 2021a). By contrast, reducing the
73 dose of disinfectant residuals in engineered water systems can mitigate the DBP issue, but
74 microbes (including pathogens) can thrive. A chlorine burn (i.e., a short conversion from
75 chloramination and free chlorination) is an effective means to control nitrification in chloraminated
76 engineered water systems (AWWA, 2013). However, chlorine burns significantly enhance DBP
77 formation in water pipes, posing serious public health risks (Alexander et al., 2024; Alfredo, 2021;
78 Allen et al., 2022).

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

89 area of research and practice ([Henrique Alves Ribeiro and Reynoso-Meza, 2023](#); [Li et al., 2021](#);
90 [Narita et al., 2023](#); [Speight et al., 2019](#)).

91 Existing reviews have summarized the applications of ML in various aspects of drinking water
92 quality ([Ewuzie et al., 2022](#); [Huang et al., 2021](#); [Zhu et al., 2022](#)), such as source water quality and
93 contamination ([Gong et al., 2023](#); [Zanoni et al., 2022](#)), the treatment processes ([Li et al., 2021](#);
94 [Lowe et al., 2022](#); [Ortiz-Lopez et al., 2022](#)), and detection of quality anomaly ([Dogo et al., 2019](#)).
95 However, a state-of-the-art, comprehensive review is missing that focuses on the applications of
96 ML to monitor, simulate, predict, and control drinking water quality, especially in engineered water
97 systems. Since engineered water systems are the vital civil infrastructure delivering municipal
98 water from water utilities to the residents and industrial/commercial consumers ([WHO, 2011](#)),
99 summarizing such applications can help understand and ensure drinking water quality, protect
100 public health, and promote socioeconomic development.

101
102 In this review, we critically compile the applications of ML in assessing and ensuring water quality.
103 We focus on engineered water systems but also cover other drinking-water-related settings such as
104 source water and drinking water treatment processes. First, we introduce ML and common ML
105 models. Then, we summarize recent progresses on the applications of ML in safe drinking water
106 supply from the physical, chemical, microbiological, and temporal perspectives. Finally, we
107 present the challenges and outlook for applying ML to ensure safe drinking waters supply. We
108 focus on the applications of ML to access and control key factors and etiologies that affect drinking
109 water quality and public health, such as water pipeline failures, DBPs, heavy metals, OPs, biofilms,
110 and antimicrobial resistance genes (ARGs). This is the first review summarizing the applications

111 of ML in assessing and controlling water quality in municipal engineered water systems, while the
112 applications of ML in other drinking-water-related settings are also discussed.

113

114 **2. Machine learning primers**

115 In the past decade, ML has driven significant progress across domains in modern society, including
116 object detection ([Erhan et al., 2014](#); [Lin et al., 2017](#)), autonomous driving ([Almalioglu et al., 2022](#);
117 [Feng et al., 2023](#)), drug delivery ([Allesoe et al., 2023](#); [Wołos et al., 2022](#)), weather forecasting ([Bi
118 et al., 2023](#)), and design-by-analogy ([Jiang et al., 2021](#)). Three key advancements in AI drive this
119 process: I) the availability of extensive datasets, II) the development of robust computing hardware,
120 and III) the refinement of advanced algorithms. In contrast to traditional physical and chemical
121 theories relying on explicit formulas for problem-solving, ML tackles problems by extracting
122 concealed insights from datasets through the learning process ([Ley et al., 2022](#)).

123

124 **2.1. Categories of machine learning**

125 On the basis of the nature of available datasets, ML can be divided into four main categories: I)
126 supervised learning, II) unsupervised learning, III) semi-supervised learning, and IV)
127 reinforcement learning ([Goodfellow et al., 2016](#)). Supervised learning uses both input data and
128 corresponding labels for training. Unsupervised learning, on the other hand, deals solely with input
129 data without labeled information. As a combination of supervised and unsupervised learning, semi-
130 supervised learning combines mostly unlabeled datasets and a few labeled ones. Reinforcement
131 learning algorithms, such as Q-learning, enable learning by interacting with an environment and
132 receiving feedback. These algorithms underpin the growth of AI, improving system performance
133 by exposing to data and experience. Among the four categories, supervised learning is the most

134 widely used and well-established in assessing and protecting drinking water quality because of its
135 strength in prediction with labeled datasets (Cordero et al., 2021; Hong et al., 2020; Zhang et al.,
136 2019c; Zhou et al., 2019). Unsupervised learning is also useful in safe drinking water supply for
137 tasks such as identifying the major factors affecting water quality and classifying data. Conversely,
138 in safe drinking water supply, the application of semi-supervised learning is scarce, and
139 reinforcement learning is unexplored.

140

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

150

151 Unsupervised learning focuses on extracting patterns, structures, and relationships from the input
152 data without labeled outputs (Goodfellow et al., 2016). Instead of targeting predefined targets,
153 unsupervised learning algorithms (such as K-means and dimensionality reduction techniques)
154 discover inherent structures within the data. For instance, K-means groups similar data on the basis
155 of their intrinsic features (Moodley and van der Haar, 2019). Dimensionality reduction techniques
156 including principal component analysis (PCA) simplify complex datasets by preserving their

157 essential characteristics (Peleato et al., 2018). In drinking water research, unsupervised learning is
158 crucial in tasks such as clustering bacteria (Moodley and van der Haar, 2019; Pinto et al., 2014),
159 simplifying data for subsequent analysis (Peleato et al., 2018), and analyzing raw data to identify
160 key parameters affecting water quality (Kazemi et al., 2023).

161

162 **2.2. Workflow chart of machine learning**

163 **2.2.1 Problem definition**

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

169

170 **2.2.2. Data collection**

171 In any ML endeavor, the quality of data is the key to the success of the modeling. Data collection
172 involves sourcing, gathering, and recording from various origins, such as observational studies,
173 controlled experiments, publications, and databases. The collected data should be pertinent to the
174 problem, accurate, and suitable for developing ML models. Along with data collection, one needs
175 to document the sources, methods, and potential biases associated with the data to ensure
176 transparency and reproducibility.

177

178 **2.2.3. Data preprocessing**

179 Data preprocessing is critical in a ML workflow that involves cleaning (i.e., filtering),
180 transforming, and organizing raw data. Cleaning is identifying and rectifying errors, inaccuracies,
181 inconsistencies, and anomalies in the collected data, which involves filtering incorrect values,
182 detecting outliers, removing duplicate entries, and converting data types. Normalization, an
183 example of data transformation, scales data into a specific range, such as 0 to 1. Normalization
184 ensures that each feature of the data contributes equally to model training, preventing any single
185 attribute from disproportionately influencing the results. Organizing raw data entails structuring
186 and optimizes the data for the specific ML algorithms to be applied.

187

188 **2.2.4. Model training**

189 The first step in model training is to divide the dataset into training, validation, and test datasets
190 (Hastie et al., 2009). The training dataset trains the model, the validation dataset tunes
191 hyperparameters, and the test dataset evaluates the model. The second step is to select an
192 appropriate model on the basis of the nature of the problem. During training, the model recognizes
193 the relationship between the inputs and outputs and minimizes the difference between predicted
194 and actual outputs by iteratively adjusting the model parameters with optimization techniques such
195 as gradient descent.

196

197 **2.2.5. Model evaluation**

198 Evaluation metrics differ on the basis of the nature of the problem (Xie et al., 2023). Evaluation
199 metrics assess the performance of classification models, revealing their ability to distinguish
200 between classes. Accuracy refers to the proportion of correctly classified instances to the total
201 instances. Precision measures the proportion of correctly predicted positive instances among all

202 predicted positives. Recall gauges the proportion of correctly predicted actual positive instances.
203 The F1-score combines precision and recall into a single metric, offering a balanced view of the
204 accuracy of a model (Sokolova and Lapalme, 2009). The confusion matrix provides a tabular
205 representation of true positive (TP), true negative (TN), false positive (FP), and false negative (FN)
206 predictions (James et al., 2013). Receiver operating characteristic (ROC) curves illustrate the
207 trade-off between the TP rate and the FP rate at different classification thresholds with the area
208 under the curve (AUC) summarizing the performance of the curve (Bradley, 1997). In regression
209 tasks, the mean squared error (MSE) and root mean squared error (RMSE) quantify the average
210 squared differences between the predicted and actual values, and the mean absolute error (MAE)
211 measures the average absolute differences (Willmott and Matsuura, 2005). Additionally, the
212 coefficient of determination (R^2) indicates the proportion of variance in the target variable
213 explained by the model (Steel and Torrie, 1960).

214

215 **3. Machine learning to ensure safe drinking water supply from the physical perspective**

216 Drinking water production and distribution is critical to public health, socioeconomic growth, and
217 urban development (Grey and Sadoff, 2007). Water demand prediction is a critical component in
218 drinking water production. Traditional methods for estimating water demand often lead to either
219 overestimation, resulting in high costs and waste of resources, or underestimation, resulting in
220 water supply shortages during peak times (Donkor et al., 2014). Additionally, water systems suffer
221 from losses because of leaks and inefficiencies with severe financial and operational consequences
222 (Lambert et al., 1999; Lee and Schwab, 2005; Reis et al., 2023). ML models are a compelling
223 solution to predict drinking water production and demand and identify system losses with greater
224 precision than conventional approaches. This section examines how ML contributes to the

225 management of water resources, the prediction of water demand, and the assurance of
226 sustainability and reliability in drinking water supply systems (Table 1).

227

228 3.1. Predicting water production and demand

229 Prediction of water production and demand is critical in safe drinking water supply. A hybrid model
230 that combines genetic algorithms (GA) and GA artificial neural networks (GA-ANN) can predict
231 drinking water production (Figure 1a) (Zhang et al., 2019c). The model uses temperature, COD,
232 and electricity and chemical consumption as the inputs. The GA-ANN was trained and validated
233 with monthly data from 45 water utilities across China. The R^2 (0.93) of GA-ANN is substantially
234 higher than that of the ANN (0.71) when more training data are incorporated. GA not only
235 optimizes the weights and biases to enhance prediction accuracy but also increases the tolerance
236 to imprecision, uncertainties, and approximates in the inputs. Moreover, GA-ANN could
237 effectively forecast fluctuations in water production for various scenarios, highlighting its
238 feasibility in adjusting water treatment operations. To assess drinking water demand patterns, one
239 can apply unsupervised learning algorithms (such as the hierarchical K-means algorithm) to raw
240 time-series data of drinking water consumption (Leitão et al., 2019). In that algorithm, daily time-
241 series demand with hourly records are the inputs within a 24-dimensional feature space to identify
242 dense and distinctly separated temporal patterns of water demand. In contrast to directly clustering
243 the water demand patterns, short-time water demand forecasting is more intriguing and has more
244 practical merits in optimizing drinking water supply. A gated recurrent unit network (GRUN)
245 predicts short-term water demand for different district metering areas for the upcoming 15 min and
246 24 h using a time-step of 15 min (Guo et al., 2018). The GRUN has a higher accuracy with a lower
247 mean absolute percentage error (MAPE) between 2.06% and 2.46% than the conventional three-

248 dense-layered ANN (MAPE between 2.46% and 2.54%) and the seasonal autoregressive integrated
249 moving average (SARIMA) model (MAPE between 2.57% and 2.85%) for the 15-min prediction.
250 For the 24-h prediction, the GRUN also achieves more precise forecasts with MAPE between 4.33%
251 and 4.96%. Another study used three ML models to forecast both short-term and long-term water
252 demand encompassing daily, weekly, and monthly intervals in Iran (Ghiassi et al., 2017). These
253 models include a dynamic artificial neural network (DAN2), a focused time-delay neural network
254 (FTDNN), and a KNN. Given its inherent design of adjusting dynamically to the data-driven
255 learning, DAN2 is promising at time-series forecasting, catering to datasets characterized by
256 evolving temporal patterns. Correspondingly, DAN2 achieves remarkable prediction accuracies
257 (96% for daily, 99% for weekly, and 98% for monthly water demand forecasts) and outperforms
258 FTDNN and KNN.

259

260 **3.2. Monitoring pipeline integrity**

261 Pipeline failures in engineered water systems cause significant water loss and contamination
262 (Renwick et al., 2019). These failures can introduce harmful microbes and chemicals from the
263 surroundings into distributed water. Addressing these issues requires precisely localizing failures
264 in complex water networks, thoroughly assessing their impacts, and preventing future incidents.
265 ML is a powerful tool to address these challenges, offering innovative solutions for detecting and
266 predicting pipeline failures. This section reviews recent advancements in the application of ML to
267 understand, detect, forecast, and mitigate pipeline failures, highlighting the role of deep learning
268 and ensemble models in this application (Table 1).

269

270 To detect pipe burst locals, a study developed a burst location identification framework by fully-
271 linear DenseNet (BLIFF) (Figure 1b) (Zhou et al., 2019), which relies on deep learning through
272 the fully-linear DenseNet (FL-DenseNet) model. BLIFF supplants the convolutional layers in
273 DenseNet with linear connections and omits pooling layers. Using real-time pressure
274 measurements as the inputs, BLIFF generates the likelihood values of a burst for each pipe in the
275 potential burst district. The prediction accuracies, ranging from 62.35% (highest probability pipe
276 match) to 98.58% (top five pipes match), of BLIFF are two times those of the original DenseNet
277 model. The remarkable improvement in the prediction accuracy is attributed to the linear-
278 connection layer, which discerns global features in the pressure signals. The effective deployment
279 of deep learning methods such as BLIFF corroborates the viability of pressure values in burst
280 localization, countering prior assertions of their insensitivity to burst events (Bakker et al., 2014;
281 Mounce et al., 2010). Another work proposed an advanced meta-learning (AdvaML) model to
282 predict the failure of water pipelines (Almheiri et al., 2021). AdvaML comprises an input layer
283 with 33 neurons (mirroring the 33 input variables including pipe and climate data), four hidden
284 layers, and an output layer that yields the failure/hazard index of a pipe. AdvaML forecasts the risk
285 index associated with pipe failures and detects pivotal determinants of pipeline service life. Of
286 these determinants, the number of traffic lanes and chlorine residual concentration are paramount,
287 collectively contributing approximately 9% to the service life analysis of water pipes. Benefiting
288 from its knowledge transfer from initial parameterization to the ultimate learning phase, AdvaML
289 has commendable performance even with scant training data compared with cox-proportional
290 hazards (Cox-PH), survival support vector machine (SSVM), and random survival forest (SRF).
291

292 While inherent system vulnerabilities cause pipeline failures, external factors exacerbate the issue
293 (Fan et al., 2023). Climatic extremes and weather disasters, such as wildfires, become more
294 frequent because of climate change and threaten drinking water infrastructures. In response,
295 researchers leverage ML to better understand and predict the impact of these disasters on water
296 pipes. Two ensemble ML models (RF and XGB) can predict the repercussions of calamities on
297 water supply infrastructures (e.g., water pipelines) (Park et al., 2020). These models incorporate
298 23 variables encompassing facility specifications and operational data from 419 water utilities in
299 South Korea. The models project the total disaster index (TDI), a metric signifying the effects or
300 damages wrought by three predominant disasters (typhoons, heavy rainfalls, and earthquakes) on
301 water supply systems. While both RF and XGB have commendable predictive prowess concerning
302 the TDI, XGB slightly outperforms in most scenarios. Another study developed four models, a
303 linear regression-based repair rate (RR) method, LR, boosted regression trees (BRT), and RF, to
304 predict pipeline damage during an earthquake (Bagriacik et al., 2018). The models incorporate
305 parameters such as ground shaking, permanent ground deformation, pipe material, pipe diameter,
306 year of installation, and trench backfill type. Each model demonstrates unique strengths. The BRT
307 model has the best overall predictive performance, while the LR model is instrumental in
308 highlighting the influence of pipe materials and trench types on pipeline damage.

309

310 **4. Machine learning to ensure safe drinking water supply from the chemical perspective**

311 Numerous chemicals, such as DBPs, disinfectant residuals, and heavy metals can appear in
312 drinking water, deteriorating water quality and affecting public health (Levallois and Villanueva,
313 2019; Valbonesi et al., 2021). Therefore, monitoring and controlling those chemicals are essential
314 to ensuring drinking water quality. Conventional approaches to assess, monitor, and/or control

315 chemicals can be time-consuming, inaccurate, and costly. ML opens a promising venue for
316 monitoring and ensuring chemical drinking water quality. This section compiles the applications
317 of ML to assess and control chemicals in drinking water with a focus on engineered water systems
318 (Table 2).

319

320 4.1. Optimizing drinking water disinfection

321 Drinking water disinfection is critical to ensuring microbial drinking water quality and
322 safeguarding public health (Zhang and Lu, 2021a). Disinfection is effective in killing pathogens,
323 impeding microbiological recontamination, and inhibiting biofilm development in drinking water
324 (Mazhar et al., 2020). Chlorine-based disinfectants, such as free chlorine (e.g., chlorine gas and
325 sodium hypochlorite), bound or combined chlorine (e.g., monochloramine), and chlorine dioxide,
326 are widely used in water treatment because of their cost-effectiveness and high efficiency (Jefri et
327 al., 2022; Zhang et al., 2018). Nonetheless, when these disinfectants interact with natural organic
328 matter (NOM) and anthropogenic compounds (such as pharmaceuticals and antimicrobials), they
329 generate DBPs such as trihalomethanes (THMs), haloacetic acids (HAAs), haloketones (HKs),
330 haloacetonitriles, halophenols, and halopropanoles (Favere et al., 2021; Xiao et al., 2023). DBPs
331 cause reproductive defects, cancer, and other serious health issues (Pandian et al., 2022; Zhou et
332 al., 2023a). Therefore, monitoring and controlling DBPs in drinking water is vital to public health
333 (He et al., 2021; Helte et al., 2023; Redondo-Hasselerharm et al., 2022). Conventional methods to
334 monitor DBPs require expensive equipment such as gas chromatography (GC) and liquid
335 chromatography (LC) combined with mass spectrometry (MS) and complicated pre-treatment
336 processes. Thus, those conventional methods are labor-intensive, costly, and time-consuming,
337 limiting the ability of water utilities to reduce DBP formation. By contrast, ML to monitor DBPs

338 in drinking water are accurate, efficient, inexpensive, and easy to handle (Table 2) (Balogun et al.,
339 2021; Jia et al., 2021; Podgorski and Berg, 2022).

340

341 4.1.1. Predicting the formation of disinfection by-products from operation conditions and 342 water quality metrics

343 Finding the optimal disinfectant dosages to minimize the levels of DBPs in finished water is crucial
344 (He et al., 2021; Zhang and Lu, 2021a). Nevertheless, reaching this goal with traditional methods
345 is time-consuming, expensive, and complex. Conversely, ML is effective in predicting the
346 formation of DBPs, significantly reducing capital and human investment for DBP control.
347 Common input parameters in these ML models are operational and water quality variables, such
348 as water temperature, contact time, pH, absorbance of light at 254 nm (UV_{254}), and the
349 concentrations of dissolved organic carbon (DOC), chloride (C_{resCl^-}), bromide (C_{Br^-}), nitrite
350 nitrogen ($C_{NO_2^-N}$), and ammonium nitrogen ($C_{NH_4^+-N}$) (Deng et al., 2021; Hong et al., 2020; Hu et
351 al., 2023; Lin et al., 2020; Pan et al., 2023; Singh and Gupta, 2012). The outputs are the
352 concentrations of DBPs, such as THMs, HAAs, and HKs.

353

354 A study developed three ML models, including ANN, SVM, and gene expression programming
355 (GEP), to forecast THM formation in chlorinated river water on the basis of a 63-point dataset
356 (Singh and Gupta, 2012). Specifically, pH, water temperature, contact time (t), bromide
357 concentration, and DOC-normalized chlorine dose (Cl_2/DOC) are the inputs. SVM outperforms
358 the other two models, exhibiting the highest R^2 and the lowest RMSE values. Furthermore,
359 sensitivity analysis reveals that pH, hydraulic retention time (HRT), and water temperature are the
360 top three contributors to DBP formation. In addition, radial basis function (RBF) based ANN

361 (RBF-ANN) can predict the formation of typical DBPs such as HAAs (Lin et al., 2020), THMs
362 (Hong et al., 2020), and HKs (Deng et al., 2021) in drinking water. For instance, a study extracted
363 64 representative data points from the literature to predict HAA formation using pH, water
364 temperature, DOC, UV₂₅₄, C_{resCl^-} , C_{Br^-} , $C_{\text{NO}_2^--\text{N}}$, and $C_{\text{NH}_4^+-\text{N}}$ as the inputs (Figure 2a) (Lin et al.,
365 2020). RBF-ANN outperforms the linear and log-linear models by 21% and 47 % in accuracy,
366 respectively. Therefore, RBF-ANN is promising in assessing DBP formation and optimizing
367 disinfection. A follow-up study used 64 data points to train an RBF-ANN to predict THM
368 formation (Hong et al., 2020). RBF-ANN achieves accuracies between 92% and 98% and
369 regression coefficients between 0.76 and 0.93, outperforming the linear and log-linear models and
370 demonstrating its superiority to uncover complex non-linear patterns in THM formation. Even
371 when trained with fewer water quality variables, a fusion of grey relation analysis with RBF-ANN
372 could provide superior prediction results. Furthermore, an RBF-ANN trained with 63 data points
373 of tap water predicts the formation of HK (Deng et al., 2021). Both RBF-ANN and back
374 propagation (BP) ANN outstrip the linear and log-linear models with the RBF-ANN displaying
375 higher accuracies in both internal and external validations. Another study applied a decision tree
376 boost (DTB) model to predict the concentrations of THM4 and HAAs (Pan et al., 2023). The study
377 correlated water quality parameters with mixed chlorine/chloramine species. The work then
378 selected seven variables such as NH_2Cl , NHCl_2 , organic chloramines, pH, total dissolved nitrogen
379 (TDN), nitrite, total organic carbon (TOC), and NH_4^+ as the independent variables to predict
380 THM4 and HAAs. The DTB model demonstrates higher prediction accuracy with R^2 values of
381 0.56 for THM4 and 0.65 for HAAs, while the inclusion of organic chloramines improves the
382 prediction. Additionally, a study implemented multiple ML models to predict emerging DBPs in
383 small water distribution systems across Canada by analyzing the data from eleven such networks

384 (Hu et al., 2023). The models use parameters such as water temperature, total chlorine residual,
385 DOC, turbidity, pH, conductivity, and UV₂₅₄ to predict the concentrations of THMs, HAAs,
386 dichloroacetonitrile (DCAN), chloropicrin (CPK), and trichloropropanone (TCP). Among the
387 evaluated models, support vector regression (SVR) and Gaussian process regression (GPR) show
388 superior performance with SVR exhibiting the highest prediction accuracy ($R^2 = 0.94$) and stability
389 for DCAN and TCP, while GPR is optimal for predicting CPK ($R^2 = 0.92$).

390

391 4.1.2. Assessing disinfection by-products using online spectroscopy

392 Fluorescence spectroscopy is the preferred technique to monitor DBPs (Krasner et al., 2006;
393 Rodriguez et al., 2004). Fluorescence spectroscopy is sensitive in assessing the characteristics and
394 reactivity of NOM because of its minimal sample preparation requirement and short acquisition
395 time (Pifer and Fairey, 2012). However, the complex high-dimensional characteristics of
396 fluorescence spectroscopy make it difficult to predict DBP formation. The significant resources
397 and time required for DBP analysis through fluorescence spectroscopy restrict the capacity of
398 water utilities to reduce DBP formation. This situation has prompted the development of ML
399 models to assess DBP formation. Compared with traditional fluorescence spectroscopy, ML-
400 powered fluorescence spectroscopy can accurately assess DBP formation with limited resource
401 and time requirement.

402

403 Autoencoder-neural networks (AE-NN) can predict the concentrations of both THMs and HAAs
404 in river water from fluorescence spectra (Figure 2b) (Peleato et al., 2018). To manage the high
405 dimensionality of the fluorescence spectra, the researchers applied three dimension-reduction
406 techniques, including AE-NN, parallel factors analysis (PARAFAC), and PCA. Afterward, they

407 trained NN to identify fluorescence regions associated with DBP formation and to predict DBP
408 concentrations. The AE-NN model has superior predictive accuracies for THMs and HAAs,
409 achieving validation MAE values of 9.65 $\mu\text{g/L}$ and 9.64 $\mu\text{g/L}$, respectively. These figures exceed
410 those of PCA, which has higher validation MAE values of 13.19 $\mu\text{g/L}$ for THMs and 11.92 $\mu\text{g/L}$
411 for HAAs. Furthermore, the precision of the AE-NN model surpasses that of PARAFAC, which
412 has validation MEA values of 20.39 $\mu\text{g/L}$ for THMs and 14.00 $\mu\text{g/L}$ for HAAs. In addition,
413 convolutional neural networks (CNN) can predict DBP concentrations from fluorescence spectra
414 without extensive data pre-processing (Peleato, 2022). Compared with multilayer perceptron
415 (MLP) and dimensionality reduction techniques, CNN not only exhibits superior prediction
416 accuracies for THMs and HAAs but also identifies the fluorescence spectra regions highly
417 associated with DBP formation.

418

419 4.1.3. Unraveling the formation mechanisms of disinfectant by-products

420 ML is promising in predicting DBP formation using either water quality and operational
421 parameters or via online spectrum monitoring. However, even with the knowledge of DBP
422 concentration, controlling DBPs in drinking water remains costly and inefficient (Bond et al., 2011;
423 Rodriguez et al., 2004). An effective approach for DBP control is to remove DBP precursors and
424 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 DBP formation.

426

427 A multiple linear regression (MLR) model can predict the production of chloroform (a THM
428 compound) from organic precursors (Bond and Graham, 2017). Relying on 211 precursors from
429 22 studies, the MLR model uses 19 descriptors as the inputs and chloroform yield as the output.

430 The well-trained MLR model has a promising prediction accuracy with an R^2 value of 0.91 and an
431 RMSE value of 8.93 mol/mol. Further chemical insights pinpoint that functional groups, such as
432 hydroxyl, chlorine, and carboxyl groups, significantly affect chloroform formation. ML can also
433 forecast the formation of HAAs from the interaction between organic precursors and free chlorine
434 (Cordero et al., 2021). The training dataset comprises 283 organic compounds and 732 chemical
435 descriptors as the inputs with HAA yield as the output. These organic compounds are converted
436 into 2D and 3D chemical descriptors with their simplified molecular input line entry system
437 (SMILE) strings used for ML compatibility. Three ML models (RF, SVR, and MLP) are selected
438 because they can handle nonlinear problems, activity cliffs, and high dimensions in addition to
439 MLR as a benchmark. RF 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 4.1.4. Evaluating alternative drinking water disinfectants

445 Since chlorine-based disinfectants produce harmful DBPs, alternative disinfectants, such as ozone,
446 for drinking water disinfection attract attention (Lin and Lin, 2024; Manasfi, 2021; Zhang et al.,
447 2019a; Zhang et al., 2019b). Unlike chlorination, ozonation does not produce chlorinated THMs
448 and HAAs. Ozone, therefore, provides a two-fold benefit: it is effective and does not generate
449 chlorinated DBPs. However, ozonation produces various other DBPs (Mao et al., 2014;
450 Richardson et al., 1999). The occurrence and toxicity of ozonated DBPs are a concern (Simpson
451 and Mitch, 2022; Srivastav et al., 2020). Therefore, while ozone does not generate chlorinated
452 DBPs, its use requires careful considerations for other toxic byproducts. This section summarize

453 how ML models have been developed and applied in the field of ozonation. ML can help control
454 the formation of bromate and reduce micropollutants, microbes, and organic contaminants during
455 ozonation.

456
457 Compared with MLR, ANN has two advantages in controlling bromate formation during ozonation
458 (Legube et al., 2004): First, ANN with an R^2 value of 0.98 is more accurate than MLR. Second,
459 ANN classifies model variables (predictors) in the descending order of impact: ozone dose, $C_{\text{NH}_4^+-\text{N}}$,
460 bromide concentration (C_{Br^-}), pH, water temperature, DOC, and alkalinity. While ANN has
461 superior performance, the simplicity of MLR is attractive. However, a key limitation of MLR is
462 that its accuracy decreases with increased sample size because MLR cannot effectively process
463 nonlinear components.

464
465 ML models based on routinely measured physical-chemical water quality parameters can predict
466 the oxidation of micropollutants (such as pharmaceuticals and personal care products) during
467 ozonation. For instance, RF can predict the oxidation of micropollutants during ozonation (Cha et
468 al., 2021) (Figure 2c). That study introduced four distinct RF models, all incorporating standard
469 predictors such as pH, alkalinity, and DOC. These models have unique inclusions of fluorescence
470 excitation-emission matrix (FEEM) data at different resolutions. These models are as FEEM-Free,
471 FEEM-LowRes, FEEM-HighRes, and FEEM-FullRes, each with a different resolution of FEEM
472 data used as the unique predictors. Integrating FEEM data results in more accurate prediction of
473 ozone exposures. The high-resolution FEEM data yield better predictions for micropollutant
474 abatement ($R^2 = 0.904$; RMSE = 6.6%). However, the improvement in prediction accuracy when

475 using FEEM data is less substantial for predicting micropollutant abatement than for predicting
476 the exposures of oxidants (i.e., ozone and hydroxyl radicals) during ozonation.

477
478 ML-quantitative-structure-property-relationship (ML-QSPR) methods can calculate the rate
479 constant (k_{O_3}) of the reactions between ozone and micropollutants (Gupta and Basant, 2016; Huang
480 et al., 2020; Shi et al., 2022; Sudhakaran and Amy, 2013). Generally, nonlinear models outperform
481 their linear counterparts. For instance, an MLR method (Sudhakaran and Amy, 2013) and an SVM
482 method (Huang et al., 2020) have R^2 values of greater than 0.75 and 0.78, respectively. Conversely,
483 a DTB model has a higher R^2 value of greater than 0.97 (Gupta and Basant, 2016). A recent study
484 compared several ML models including MLR, SVM, DT, RF, and deep neural network (DNN) for
485 predicting $\log k_{O_3}$ (Shi et al., 2022). Of these, RF has the highest effectiveness with a peak R^2 value
486 of 0.91. RF has two primary benefits: robustness and a lower tendency of overfitting. On the other
487 hand, DT has a complex structure and subsequently increases the overfitting risk. In addition, DNN,
488 promising at recognizing nonlinear features, underperforms in predicting $\log k_{O_3}$. A similar
489 situation occurs when ML models predict the elimination of recalcitrant trace organic compounds
490 (TOrcs) by ozonation for municipal wastewater reuse (Park et al., 2015). Specifically, ANN is
491 susceptible to overfitting. Incorporating PCA into ANN creates a PC-ANN workflow, which
492 addresses the overfitting issue. PCA transforms the input variables (Table 2) to linearly
493 independent variables, thereby resolving the issue of collinearity among explanatory variables.
494 The PC-ANN workflow ($R^2 = 0.934$) surpasses the standalone ANN ($R^2 = 0.914$) regarding
495 predictive power.

496

497 4.2. Surveilling drinking water nitrification

498 Chloramine, a commonly used chlorine-based disinfectant, can maintain a higher level of residual
499 while minimizing DBP formation (Shao et al., 2023; Shi et al., 2020). However, nitrification is a
500 major concern in chloraminated engineered water systems (Allen et al., 2022; AWWA, 2013).
501 During nitrification, ammonia-oxidizing microbes oxidize free ammonia to nitrite, and nitrite-
502 oxidizing bacteria further oxidize nitrite to nitrate. Nitrification deteriorates water quality by
503 destroying chloramine residuals, releasing free ammonia, promoting microbial (re)growth, and
504 producing toxic nitrite and nitrate. Therefore, monitoring and controlling nitrification in
505 engineered water systems is critical to ensuring drinking water quality and protecting public health.
506
507 ML is useful in detecting drinking water nitrification (Table 2). NB classifier, a supervised ML
508 model, relies on biomass and microbiome datasets to detect nitrification in engineered water
509 systems (Gomez-Alvarez and Revetta, 2020). After being trained with microbial indicators, the
510 model has a binary classification accuracy of up to 85% with an AUC of 0.825 when distinguishing
511 between nitrification and stable events. Nitrification can also be monitored using spectrum
512 fingerprint since one can isolate the combined nitrate and nitrite from the total spectra. SVR was
513 used and trained to predict the concentrations of nitrate and nitrite from nitrate/nitrite spectra at
514 various wavelengths (Figure 2d) (Hossain et al., 2021). SVR negates the need for any chemical
515 supplements, is easy to use, and can reach a high level of precision of up to ± 0.01 mg N/L.

516

517 **4.3. Monitoring and regulating heavy metals in drinking water**

518 The detachment of heavy metals from water pipes (i.e., leaching) deteriorates drinking water
519 quality in engineered water systems (Mays, 2000; Proctor et al., 2020). Heavy metals are toxic to
520 human beings and significantly affect public health (Abd Elnabi et al., 2023; Fu and Xi, 2020).

521 ML is useful in assessing heavy metals in drinking water (Yaseen, 2021; Zhu et al., 2022).
522 Therefore, this section summarizes the applications of ML in monitoring and regulating heavy
523 metals in drinking water (Table 2).

524

525 4.3.1. Assessing heavy metal concentration and distribution

526 A continuous on-site, *in situ* system can estimate lead (Pb) concentration in municipal water (Oh
527 et al., 2021). The system leverages the SVR algorithm, supplanting traditional mathematical
528 models confined to analyzing stationary ions in a solid substrate. By using the radio-frequency
529 reflection coefficient of the raw trace data, the system predicts Pb ion concentration with a
530 resolution of 1 µg Pb/L and an RMS prediction error of 0.71 µg Pb/L in the presence of interfering
531 metals such as copper (Cu²⁺), ferric (Fe³⁺), and zinc (Zn²⁺) ions. Other than estimating heavy metal
532 concentration in individual samples, ML is promising in broader analytical applications. For
533 instance, ML is useful in spatially interpolating environmental variables, significantly enhancing
534 its performance (Li et al., 2011). This approach is valid in developing spatial interpolation maps
535 depicting the concentrations of heavy metals such as Fe, Mn, Ni, Pb, and Zn in groundwater (i.e.,
536 a source of drinking water) (De Jesus et al., 2021). That study combined ML and geostatistical
537 interpolation (MLGI) to leverage an ANN-based algorithm to augment the efficacy and robustness
538 of the spatial interpolation mapping. Furthermore, the MLGI approach comprehensively assesses
539 the carcinogenic risks of heavy metals through *in situ* measurements. The approach produces
540 detailed spatial maps delineating heavy metal concentration and estimates health quotient indices
541 (HQI) to offer a more refined risk assessment (Senoro et al., 2022). While integrating ML
542 algorithms elevates the efficacy and robustness of spatial interpolation, traditional interpolation
543 techniques are still critical in this domain. For instance, a spherical semi-variogram model relying

544 on the classic Kriging interpolation technique can monitor the temporospatial distribution of
545 residual aluminum (Al) in a DWDS, highlighting the enduring relevance and applicability of
546 traditional interpolation methods (Tian et al., 2020).

547

548 4.3.2. Predicting heavy metal removal with porous materials

549 Adsorption is proficient in mitigating heavy metal contamination in drinking water (Joseph et al.,
550 2019; Wołowiec et al., 2019). The adsorption of heavy metals by porous materials has highly
551 stochastic, non-linear, and non-stationary dynamics coupled with redundancy (Bhagat et al., 2020).
552 Therefore, ML is a preferred technique for analyzing the removal of heavy metals of porous media.
553 Many ML techniques can enhance the precision and efficacy of predicting heavy metal adsorption
554 dynamics by porous materials. Common predictors for those ML techniques are adsorbent dosage,
555 operating temperature, contact time, and pH, whereas the output is the removal of heavy metals.
556 Other variables can also be incorporated into the predictive models such as the initial concentration
557 of heavy metals, the specific surface area of metal-organic frameworks (MOFs), and the presence
558 of anions. A study used four tree-based ML models, including light gradient-boosting machine
559 (LightGBM), XGB, gradient-boosted decision trees (GBDT), and RF, to predict the adsorption of
560 arsenate [As(V)] by MOFs (Abdi and Mazloom, 2022). Among these models, LightGBM yields
561 the most accurate and reliable prediction with R^2 and RMSE values of 0.996 and 2.069,
562 respectively. The sensitivity analysis indicates that the adsorption process is adversely affected by
563 the initial As(V) concentration and is directly influenced by the specific surface area and dosage
564 of MOFs. ANN and symbiotic organisms search (SOS) algorithm can predict the removal of five
565 heavy metals (Al, Cd, Co, Cu, Fe, and Pb) by two adsorbents, Chitosan and Chitosan-
566 Montmorillonite nanocomposite (Hamidian et al., 2019). First, RBF-ANN has a higher prediction

567 accuracy than MLP-ANN for the two adsorbents and the five heavy metals. Second, when
568 integrated with SOS algorithm, RBF-ANN facilitates the identification of optimal performance
569 parameters and increases the adsorption performance than the experimental results. Finally, RBF-
570 ANN outperforms Langmuir and Freundlich models when considering the five heavy metals and
571 three operational parameters (pH, absorbent dose, and contact time).

572 **5. Machine learning to ensure safe drinking water supply from the microbiological** 573 **perspective**

574 Microbial drinking after quality is essential to public health and the economic development of
575 society (Abkar et al., 2024; Figueras and Borrego, 2010; Wen et al., 2020). The microbial
576 community in drinking water is highly diverse and ever-changing spatiotemporally, especially in
577 engineered water systems (Ashbolt, 2015; Jing et al., 2023; Zhou et al., 2023b). Therefore,
578 monitoring and ensuring microbial drinking water quality is complex, time-consuming, and often
579 ineffective with conventional assays. Microbes in drinking water encompass general heterotrophic
580 bacteria, fecal contaminants, microbial indicators, protozoans (such as amoebae, ciliates, and slime
581 molds), and OPs (Abkar et al., 2024; Siponen et al., 2024; Zhang et al., 2021b). Because of the
582 complex nature of microbial drinking water quality, ML models are the preferred approaches for
583 analyzing and ensuring microbial drinking water quality (Mahajna et al., 2022; Naloufi et al., 2021;
584 Saboe et al., 2021; Zhu et al., 2022). This section summarizes the applications of ML in monitoring,
585 predicting, and ensuring microbial drinking water quality (Table 3).

586

587 **5.1. Surveilling and mitigating opportunistic pathogens**

588 In municipal water, OPs are the most significant aspect of microbial drinking water quality because
589 of their frequent occurrence, high concentration, high resistance to disinfectant residuals, and

590 proliferation within amoebae (Zhang and Lu, 2021b). OPs are the major disease-causing agents in
591 drinking water and significantly affect the health of the end consumers. Therefore, closely
592 monitoring OPs in drinking water, especially in engineered water systems, is critical to assessing
593 drinking water quality and protecting public health. Dominant water-related OPs are *Legionella*
594 (especially *L. pneumophila*), *Mycobacterium* (e.g., nontuberculosis mycobacteria or NTM and *M.*
595 *avium* complex or MAC), *Pseudomonas aeruginosa*, *Vermamoeba vermiformis*, *Naegleria fowleri*,
596 and *Acanthamoeba* (Donohue et al., 2019; Isaac and Sherchan, 2020; Lytle et al., 2021). Among
597 them, *Legionella* is the most important OP. In addition, compared with conventional microbial
598 drinking water quality indicators such as fecal coliforms and , *Legionella* is a better candidate to
599 indicate microbial drinking water quality (Zhang and Lu, 2021b). Therefore, this section focuses
600 on the applications of ML in monitoring and controlling *Legionella* in municipal water (Table 3).
601
602 Studies using ML to assess the risks of OPs in drinking water remain limited. An early work
603 mitigated the proliferation of *Legionella* in premise plumbing by controlling environmental
604 variables (Sincak et al., 2014). Using water flow and water temperature as the inputs, that study
605 presented a NN-based simulator relying an approximate reasoning architecture (NARA) neuro-
606 fuzzy system to predict and simulate water tank temperature. The simulator emulates conditions
607 that inhibit the spread of *Legionella* in water networks. The NARA-based simulator achieves a
608 high fidelity in mimicking water tank temperature with an accuracy exceeding 97%. A recent study
609 integrated both unsupervised and supervised ML to correlate the spread of *Legionella* with
610 environmental variables in retirement homes, health-related facilities, tourism-related buildings,
611 and swimming-pools s from 2002 to 2019 in Italy (Brunello et al., 2022). That study used an
612 unsupervised ML algorithm to identify the spatiotemporal distribution of atypical *Legionella*

613 through an ordinal regression model. The results indicate how the distribution is correlated with
614 the types of healthcare facilities. The propagation of *Legionella* and both the nature of the facilities
615 and broader geographical characteristics have strong correlations. Hospitals have the highest
616 contamination cluster locations That work also used supervised ML to assess the serotypes of
617 *Legionella* and to anticipate the corresponding contamination levels. For serogroup assessment,
618 XGBoost, LR, and SVM Classifier were used and compared. XGBoost shows superior
619 performance with an overall classification accuracy of 0.71. The Shapley values evaluates the
620 contribution of each predictor to the final classification. The Shapley values quantify the
621 contribution of each variable to the outputs of a ML model by comparing the effect of the outputs
622 relative to the average across all inputs. The geographical location of a sample is the most
623 important parameter but is useful only when combined with other predictors. For contamination
624 level prediction, all three models demonstrate low performance with the highest accuracy of 0.57
625 from XGBoost.

626

627 **5.2. Analyzing drinking water microbial communities**

628 The microbial community in drinking water is complex with ever-changing structure and
629 significant public health implications (Abkar et al., 2024; Zhang et al., 2021a). Assessing the
630 structure and composition of the microbial community helps understand and ensure microbial
631 drinking water quality. Conventional approaches can monitor the microbial community in drinking
632 water, such as phenotypic/genotypic matching, molecular marking, high-throughput metagenomic
633 sequencing, and microbial and chemical indication. However, these methods have limitations in
634 terms of cost, time, and spatial/temporal coverage. ML, on the other hand, can overcome those

635 limitations and are suitable for analyzing, monitoring, and source-tracking microbial communities
636 in municipal water (Table 3).

637
638 When using the NB theorem to estimate the distribution of microbes in water sources, one could
639 apply either maximum posterior probability (evaluation metrics: $RMSE_c$) (Ritter et al., 2003) or
640 direct averaging posterior probability (evaluation metrics: $RMSE_p$) (Greenberg et al., 2010). Direct
641 averaging of the source posterior probability yields more precise source distribution estimates with
642 $RMSE_c$ being significantly lower than $RMSE_p$. The more precise source distribution estimates are
643 because direct estimation bypasses the information loss that typically happens when frequencies
644 are first classified and then averaged (Greenberg et al., 2010). SourceTracker as a ML tool
645 estimates the proportion of contaminants (Knights et al., 2011). SourceTracker employs the Gibbs
646 sampling technique within a Bayesian framework and is more efficient than both NB- and RF-
647 based source tracking methods (Smith et al., 2010). The superior performance of SourceTracker is
648 because it can handle ambiguity in the source and sink distributions and can model a sink sample
649 as a blend of various sources. SourceTracker can track the origin of bacteria in drinking water and
650 water sources (Liu et al., 2018). For instance, a study developed six ML models (XGBoost, KNN,
651 NB, SVM, NN, and RF) to predict microbial contamination in a watershed in California using data
652 on land cover, weather, and hydrologic variables (Wu et al., 2020). That study used SourceTracker
653 to generate ground-truth data for training purposes. Out of the six models, XGBoost outperforms
654 the other models in terms of accuracy and AUC (average AUC = 0.88) when tracking the primary
655 sources of microbial contamination in the watershed.

656

657 Unsupervised ML can unveil hidden features, trends, or patterns in bacterial communities in
658 engineered water systems. For instance, alpha and beta diversity analyses can display the spatial
659 dynamics and temporal trends of bacterial communities in DWDSs (Pinto et al., 2014). UniFrac
660 as an unsupervised ML tool uses principal coordinates analysis (PCoA) coupled beta diversity
661 measure to analyze the differences among microbial communities (Lozupone et al., 2011). UniFrac
662 can effectively analyze the microbiome in drinking water (Bruno et al., 2018; Li et al., 2017; Ling
663 et al., 2018).

664

665 5.3. Detecting drinking water parasites

666 *Cryptosporidium* and *Giardia* are protozoan parasites in municipal water with substantial public
667 health risks by causing cryptosporidiosis and giardiasis, respectively (CDC, 2021a, 2021b). These
668 pathogens are highly resistant to disinfectants, challenging drinking water treatment (Adeyemo et
669 al., 2019). Therefore, detecting and controlling *Cryptosporidium* and *Giardia* is critical to
670 maintaining drinking water quality. In this section, we discuss the performance of ML models in
671 monitoring *Cryptosporidium* and *Giardia* in drinking water (Table 3).

672

673 ML to detect *Cryptosporidium* and *Giardia* are robust and precise. For instance, deep-learning-
674 based image classification models such as ParasNet (Xu et al., 2020) and MCellNet (Luo et al.,
675 2021) are accurate in detecting these two parasites in drinking water. They have the power of ML
676 in classifying parasites from the cell-level scattering images. In addition, a linear ML model can
677 predict the contamination of these two parasites in surface water and drinking water (Ligda et al.,
678 2020), offering a valuable tool to control waterborne diseases.

679

680 ParasNet (Xu et al., 2020) uses an eight-layer CNN to determine whether particles in cell-level
681 scattering images from drinking water are *Cryptosporidium* and *Giardia*. The model has superior
682 performance compared with a traditional handcraft SVM regarding both detection accuracy and
683 processing time. For instance, ParasNet can reach a detection accuracy of above 95.6% with
684 analysis speeds of up to 100 frame-per-second (fps) on embedded Jetson TX2 platform. MCellNet
685 (Luo et al., 2021), another image classification pipeline, uses a DNN optimized from MobileNetV2
686 (Sandler et al., 2018) to recognize objects. MCellNet includes a convolutional layer, six inverted
687 residual blocks (IRBs), a flattened layer, and a fully connected layer. MCellNet can process images
688 from flow cytometry to classify *Cryptosporidium* and *Giardia*. Compared with ParasNet,
689 MCellNet achieves a higher detection accuracy of above 99.6% with a 346-fps analysis speed. The
690 superior accuracy and fast analysis of MCellNet are due to the cascading six IRBs.

691
692 An alternative statistical model uses linear discriminant function analysis (LDFA) to predict the
693 appearance of *Cryptosporidium* and *Giardia* in drinking water (Ligda et al., 2020). That model
694 uses microbiological, physicochemical, and meteorological parameters to classify the
695 contamination of *Cryptosporidium* and *Giardia* into four categories: none, low, moderate, and high
696 (oo)cysts concentrations. LDFA has accuracies of 75% and 69% in predicting *Cryptosporidium*
697 and *Giardia*, respectively.

698

699 **5.4. Assessing biofilm development in engineered water systems**

700 Biofilms are complex microbial communities adhering to surfaces (Flemming and Wingender,
701 2010). In engineered water systems, biofilms develop on the inner wall of water pipes and pose
702 significant risks. Drinking water biofilms harbor pathogenic and antimicrobial resistant bacteria,

703 corrode pipes, reduce water flow rate, deteriorate water quality, and increase the costs and
704 complexity of water distribution (Simoes and Simões, 2013; Wingender and Flemming, 2011;
705 Zhou et al., 2023c). Monitoring and controlling biofilm formation in engineered water systems
706 with conventional approaches are challenging because assessing the biofilms inside of the water
707 pipes is difficult and biofilms are protected by an extracellular matrix (Flemming et al., 2023;
708 Karygianni et al., 2020). ML provides novel solutions for assessing and controlling drinking water
709 biofilms (Table 3).

710
711 ML models studying biofilm development in engineered water systems use relevant physical (such
712 as water age, flow velocity, hydraulic regime, pipe material, and pipe age) factors to assess the
713 dynamics of biofilm development, where the heterotrophic plate count (HPC) is the output.
714 Established ML algorithms are preferred models to study the dynamics of biofilm development,
715 such as NB, RT, and RF (Ramos-Martínez et al., 2014, 2016). These algorithms have high
716 prediction accuracy and provide an in-depth understanding of the impact of physical factors on
717 biofilm development in engineered water systems. For instance, a Bagging naïve Bayesian tree
718 (B-NBT) model proposes optimal flow velocities for different types of pipes to mitigate biofilm
719 development (Ramos-Martínez et al., 2014). Predicted biofilm development probabilities show
720 that, to control biofilm accumulation, water utilities need to avoid cement pipes, implement
721 medium- or high-flow velocities in metal pipes, and sustain water ages above 0.035 in plastic pipes.
722 The 'water age' is a synthetic index derived from the normalized HRT and the distance from the
723 disinfection source.

724

725 Recent studies have enabled more detailed, single-cell level analyses and predictions of biofilm
726 development in engineered water systems (Berne et al., 2018). In addition, studies have expanded
727 ML algorithms to incorporate deep learning (Jelli et al., 2023; Weigert et al., 2020). These
728 innovative approaches can enhance the assessment of biofilm dynamics in various settings
729 including engineered water systems. A recent work (Jelli et al., 2023) optimized StarDist (Weigert
730 et al., 2020), a cutting-edge CNN-based segmentation algorithm, to segment single cells in
731 biofilms, track cell lineages, and measure single-cell growth rates (Figure 3a). First, an iterative
732 semi-automated annotation workflow was developed to accelerate the annotation of bacterial cells
733 in 3D images for training data. Then, a new post-processing algorithm (StarDist OPP) that
734 reconstructs the bacterial cell shapes were developed to increase the accuracy of bacterial
735 segmentation. The second step was to overcome the embedded limitations of the StarDist non-
736 maximum-suppression post-processing that considers only the shape information in the voxel with
737 the highest assigned label probability. StarDist OPP achieves unprecedented accuracy in biofilm
738 segmentation, surpassing other algorithms under scrutiny, such as Cellpose (Stringer et al., 2021),
739 a multi-class U-Net (Zhang et al., 2020), and BCM3D 2.0 (Zhang et al., 2022). Finally, the accurate
740 single-cell segmentation results were used to track cell lineages and to spatiotemporally measure
741 single-cell growth rates.

742

743 **5.5. Analyzing the risks and tracking the sources of antimicrobial resistance**

744 Antimicrobials have been extensively used since the 1920s in the medical industry, animal
745 husbandry, consumer goods production, and other fields (Chang et al., 2015; Hutchings et al., 2019;
746 Prescott, 2017; WHO, 2021). However, a significant portion of antimicrobials consumed by
747 humans and animals is not metabolized but excreted, entering waterbodies. Antimicrobials in

748 waterbodies contributes to the development of antibiotic-resistant bacteria (ARB) and ARGs,
749 posing a serious threat to aquatic ecosystems and public health by causing antimicrobial resistance
750 (AMR) (Roca et al., 2015; Walesch et al., 2023). The advent of ML has introduced novel
751 methodologies that enhance our ability to assess the risks and track the sources of AMR in in
752 drinking water with unprecedented precision (Table 3).

753

754 **5.5.1. Assessing the risks of antimicrobial resistance**

755 AMR presents challenges in treatment, costs, and mortality rates compared with non-resistant
756 infections in humans and animals. Despite the efforts to assess and control AMR, challenges
757 remain because of uncertainties in data acquisition and dose-response mechanisms. To streamline
758 the estimation process and minimize labor, a study developed three ML models (LR, DT, and RF)
759 to rapidly predict the relative risks of AMR in drinking water (Wu et al., 2022). These models take
760 four land-use type factors (residential, urban, green, and agriculture) and eleven environmental
761 factors (water temperature, pH, oxidation-reduction potential, electrical conductance, resistivity,
762 total dissolved solids, salinity, pressure, DO, turbidity, and 24-h accumulated rainfall) as the inputs.
763 Given limitations in data, particularly in field data, employing classification over direct regression
764 for relative risk assessment is more robust. That study used a binary classification framework,
765 labeling relative risk scores above the median as 1 for relatively high risks and scores below the
766 median as 0 for relatively low risks. Compared with LR and DT, RF has the highest accuracy
767 (0.86), precision (1.0), recall (0.75), F1 (0.86), and AUC (0.88). Finally, the feature importance
768 analysis from RF reveals that green (areas designated for natural vegetation, parks, forests, or other
769 green spaces), DO, and pH are the top-three significant influencing factors of AMR in drinking
770 water.

771

772 **5.5.2. Tracking the sources of antimicrobial resistance**

773 Leveraging the capabilities of ML for environmental monitoring, recent studies have harnessed
774 SourceTracker, an ML tool based on a Bayesian classification algorithm (Knights et al., 2011), to
775 identify the sources of ARGs in drinking water and water sources. A study used SourceTracker to
776 identify the complex sources of ARGs and assessed their contributions to ARG pollution in a peri-
777 urban river (Chen et al., 2019). The results show that the discharge from sewage plants was the
778 largest contributor of ARGs (81.6% to 92.1%) in the river sediments. Another work used
779 SourceTracker to monitor the presence of ARGs in household drinking water and tracked their
780 origins back to anthropogenic sources, highlighting the significant impact of human activities on
781 drinking water quality (Figure 3b) (Wang et al., 2023). The data generated by SourceTracker have
782 a strong Pearson correlation ($r = 0.98$) with the corresponding expected proportion by artificial
783 source inputs. Source tracking analysis from that study indicates that a significant proportion of
784 ARGs (37.1%) was from anthropogenic sources, especially wastewater effluent.

785

786 **6. Machine learning to ensure safe drinking water supply from the temporal perspective**

787 The increasing use of ML in safeguarding drinking water quality has led to the development of
788 innovative approaches to detect drinking water quality from the temporal perspective (Table 4)
789 (Zhong et al., 2021). The term 'temporal' refers to the time-related applications of ML to track,
790 predict, and mitigate contamination events in drinking water as they unfold over time. By
791 examining the temporal patterns and trends of contamination events, we can bolster the predictive
792 power and responsiveness of ML to ensure effective measures against accidental contamination in
793 drinking water. This section explores the advancements, versatility, and potential of ML in

794 revealing drinking water quality from the temporal perspective focusing on accidental drinking
795 water contamination events.

796

797 Three studies developed multiple ML approaches to detect anomalies in the drinking water quality
798 datasets from GECCO Industrial Challenges (GECCO IC) (Fehst et al., 2018; Muharemi et al.,
799 2019; Qian et al., 2020). These studies pinpointed shifts or spotted anomalies in drinking water
800 quality over time. Various parameters such as pH, redox potential, electric conductivity, turbidity,
801 and chlorine dioxide concentration are the predictors, whereas events in Boolean form are the
802 outputs. A study developed SVM, DNN, long short-term memory (LSTM), recurrent NN (RNN),
803 LR, simple NN, and linear discriminant analysis (LDA) to detect water quality anomaly in the
804 dataset from 2017 GECCO IC (Muharemi et al., 2019). SVM shows the highest performance with
805 an F1-score of 0.99 in cross-validation. Nevertheless, all the models have poor performance with
806 the unseen test dataset with a maximum F1-score of 0.36. In the other two studies focusing on the
807 dataset from GECCO IC 2018, LSTM demonstrates superior results, scoring a higher F1-score
808 than traditional models such as LR and SVM with F1-scores of 0.80 and 0.78, respectively (Fehst
809 et al., 2018; Qian et al., 2020).

810

811 The existing research, including the three key studies using the GECCO IC datasets (Fehst et al.,
812 2018; Muharemi et al., 2019; Qian et al., 2020), has made significant progress in understanding
813 anomaly in drinking water through ML. A notable trend in recent research is using real-time or
814 online applications to reflect a crucial evolution toward practical, real-world implementations.
815 Specifically, a study implemented an LSTM-based approach to detect anomalies in water quality
816 focusing on turbidity and conductivity (Rodriguez-Perez et al., 2020) (Figure 4a). That study

817 highlights the efficacy of semi-supervised classification, which retains only normal values, in
818 identifying abrupt changes and minor spikes in water quality. By contrast, supervised classification,
819 which considers both normal and anomalous data, is more suitable in identifying long-term
820 anomalies linked to gradual changes. Notably, the LSTM-based approach surpasses regression-
821 based autoregressive integrated moving average (ARIMA) in detecting these long-term anomalies.
822 Another study introduced an innovative stacking ensemble model designed for contamination
823 detection (Figure 4b) (Li et al., 2022). The model uses various water quality parameters such as
824 total chlorine, pH, electrical conductivity, water temperature, TOC, and turbidity. That model
825 integrates multiple predictors into a meta-predictor, trained through cross-validation. That
826 approach enhances the ability of the model to discern distinct features across water quality
827 parameters. The ensemble has predictors such as ANN, SVM with a linear kernel, linear regressor,
828 extra trees, uniform weighted KNN, and an RF meta-predictor. The ensemble demonstrates
829 superior performance in detecting contamination compared with an ANN benchmark method,
830 achieving higher accuracy, lower false positive rates, and improved F1-scores.

831
832 However, these models focus on single-site, one-dimensional time series data, neglecting the
833 spatial relationships inherent in multi-site sensor data. This limitation could increase false alarm
834 rates, particularly under conditions of high hydraulic variability. To address this issue, a follow-up
835 study proposed a novel unsupervised, generative-adversarial-networks-based (GAN-based)
836 multivariate method to detect multi-site contamination events (Figure 4c) (Li et al., 2023). That
837 method effectively captures spatiotemporal patterns by transforming water quality data from single
838 and multiple sites into superimposed images. The GAN-based model, having a generator and a
839 discriminator, evaluates the degree of abnormality at each time step by generating anomaly scores.

840 The generator is trained to map historical image data to expected current images, while the
841 discriminator differentiates between generated and actual normal images. That method is
842 benchmarked against a multivariate unsupervised method using a minimum-volume-ellipsoid
843 (MVE)-based event detection model (Oliker and Ostfeld, 2014). That method demonstrates
844 superior performance in all contamination scenarios, including enhanced detection rates and
845 reduced false alarms, particularly for sensor groups positioned at varying distances from the
846 contamination source. Another unique ML approach can rapidly signal potential contamination
847 risks in drinking water (Asheri Arnon et al., 2019). That approach uses an algorithm for the early
848 detection of drinking water contamination against an unpredictable stochastic background. By
849 extracting key features from the spectrophotometric characteristics of water, the algorithm can
850 effectively identify contamination using a unique affinity measure (Asheri-Arnon et al., 2018). The
851 measure compares the absorbance spectra of different water sources, thereby amplifying the
852 feature dissimilarity between portable and contaminated water, followed by processing via SVM
853 and post-processing. That chain of data processing generates a reliable early warning system for
854 contamination events with low false positives and high true alarm accuracy. The pre-processing
855 stage (the affinity measure and amplification) is essential to achieving high accuracy but may be
856 unnecessary to obtain minimal false positives.

857

858 **7. Machine learning model distribution in safe drinking water supply**

859 We provide a macroscopic visual illustration to elucidate the distribution of ML models across
860 research topics in safe drinking water supply (Figure 5). To facilitate a clear and concise visual
861 representation, we group certain ML models under broader principal categories on the basis of
862 their foundational architecture. For instance, models such as GA-ANN, Multi-layered-ANN, and

863 DNN share foundational characteristics inherent with ANN. Consequently, to elucidate the
864 overarching trends in model preferences across studies, we categorize these models as “NN-based.”
865 This approach discerns the broader trends and preferences in ML model selection and also
866 highlights the potential commonalities across research endeavors.

867

868 NN-based and regression-based ML models are the top two frequently implemented in safe
869 drinking water supply. NN-based models have significant applications in managing the production
870 and demand of drinking water and accessing and controlling DBPs. The prominent role of NN-
871 based models in these two fields is not coincidental but rather from the synergy between the
872 inherent characteristics of these fields and the strengths of NN-based models. Water management
873 and DBP assessment often involve multifaceted, nonlinear, and high-dimensional data that demand
874 robust modeling ([Aliashrafi et al., 2021](#); [Ates et al., 2022](#); [Ghobadi and Kang, 2023](#)). Given their
875 capability to model complex non-linear relationships and handle various intricate data, NN-based
876 models are an optimal solution in these contexts. For instance, the unpredictability and variability
877 in water demand patterns or the multifarious factors influencing DBP formation both require a
878 model that can discern patterns from large, intricate datasets ([Ahmadpour et al., 2023](#); [Avni et al.,](#)
879 [2015](#)). Furthermore, the flexibility of NN-based models in accommodating changing inputs makes
880 them promising in assessing the dynamic nature of drinking water quality. The wide applications
881 of NN-based models in safe drinking water supply are due to this harmonious fit between the
882 challenges posed by these fields and the advantages of these models.

883

884 By contrast, while regression-based models are widely applied in drinking water research, they
885 have suboptimal performances in certain contexts ([Almheiri et al., 2021](#); [Deng et al., 2021](#); [Hong](#)

886 et al., 2020; Legube et al., 2004; Lin et al., 2020; Rodriguez-Perez et al., 2020). The suboptimal
887 performances do not undermine the value of regression-based models. However, their linear or
888 predefined non-linear structures may limit their effectiveness, especially when compared with the
889 adaptive and intricate abilities of NN-based models.

890

891 The superior performance of NN-based models is widely acknowledged (Goodfellow et al., 2016).

892 These general strengths become pertinent when NN-based models are applied to the complexities

893 of drinking water research. First, unlike regression-based models which are limited by their linear

894 or defined non-linear structures, NN-based models capture intricate, non-linear associations.

895 Second, the mutable architecture of NN-based models allows them to modify their framework

896 during training, optimizing alignment with the inherent data distribution. Lastly, given abundant

897 data, NN-based models excel in discerning subtle data patterns because of their proficiency in

898 processing high-dimensional input attributes, whereas regression-based models may have

899 tendencies of underfitting. This proficiency of NN-based models is further enhanced by the use of

900 techniques such as grid search for hyperparameter optimization, particularly crucial in fine-tuning

901 the performance of NN-based models because of their complex architectures and the numerous

902 parameters required (Daniel et al., 2023; Rodriguez-Perez et al., 2020).

903

904 CNN-based models represent a specialized subclass of NN-based models adept at discerning

905 patterns in images or other forms of multi-dimensional data (LeCun et al., 1989; Lecun et al., 1998).

906 Therefore, we list the CNN-based models out of the broader NN category (Figure 5). The practical

907 implications of CNN-based models are evident in drinking water research: They can interpret 2D

908 fluorescence spectra and predict the formation of DBPs during disinfection (Peleato, 2022),

909 classify microbes using cell-level scattering images from drinking water (Luo et al., 2021; Xu et
910 al., 2020), and identify cells in 3D drinking water biofilm images (Jelli et al., 2023).
911

912 Other ensemble approaches are also widely applied in safe drinking water supply such as RF
913 (Breiman, 2001), XGB (Chen and Guestrin, 2016), boosted decision trees (BDT) (Friedman, 2001),
914 and stacking model (Wolpert, 1992). The core strength of these ensemble techniques is their ability
915 to amalgamate predictions from several models, aiming to boost accuracy and diminish overfitting
916 (Hastie et al., 2009). In drinking water research where data can be noisy, varied, and sometimes
917 sparse, such strategies are invaluable. Several comparative studies have delved into the
918 performance nuances of different ensemble models. A recurring observation in these investigations
919 is that the slight edge XGB outperforms RF (Abdi and Mazloom, 2022; Park et al., 2020; Wu et
920 al., 2020). Furthermore, BRT outperforms RF (Bagriacik et al., 2018). Interestingly, while XGB
921 has consistent prowess, LightGBM, another gradient boosting framework, outperforms XGB
922 (Abdi and Mazloom, 2022). Therefore, as gradient boosting algorithms continue to evolve, newer
923 iterations such as LightGBM offers even more refined performance. However, while ensemble
924 methods offer certain advantages, their efficacy is not universally dominant across scenarios. The
925 best model is often contingent upon the nature of the problem, the characteristics of the data, and
926 the specific objectives of the study. Ensemble models, with their ability to amalgamate insights
927 from multiple “weak learners,” might excel in scenarios where data are diverse, noisy, and/or
928 sparse (Fasel et al., 2022; Pang et al., 2018; Sluban and Lavrač, 2015). By contrast, for problems
929 where the data structures are deeply hierarchical or when data patterns are straightforward, NN-
930 based models or regression-based models are more suitable. The crucial factor is to match the
931 ability of the models with the specific demands and characteristics of the data sets.

932

933 We include (S)ARIMA, Kriging interpolation, SaTScan, LDFA, alpha and beta diversity analyses,
934 UniFrac, and MVE in statistical models (Figure 5). These models are more deterministic and often
935 rooted in foundational principles and established, theoretical, and/or empirical relationships. For
936 instance, SARIMA and Kriging interpolation can capture temporal and spatial patterns,
937 respectively (Guo et al., 2018; Tian et al., 2020). Alpha and beta diversity analyses and UniFrac
938 quantify microbial community diversity and compositional differences (Li and Wu, 2019). These
939 models typically operate under specific assumptions about the underlying data distribution or
940 spatiotemporal relationships. By contrast, ML, especially deep learning, is more adaptive, learning
941 patterns directly from the data without stringent assumptions (Khattak et al., 2022; Savadatti et al.,
942 2022; Singh et al., 2023).

943

944 **8. Challenges and outlooks**

945 While ML has made significant progress in drinking water research, several areas remain untapped,
946 offering significant potential for exploration and improvement. Crucial topics, such as biofilm
947 development, the assessment of AMR risks, and the evaluation of pathogen-related dangers in
948 engineered water systems, are not fully explored. The untapped potential in these fields is immense,
949 and the need to bridge the interdisciplinary divide is critical.

950

951 One significant barrier is the disconnect between water experts and AI specialists. Water scientists
952 and engineers may not be conversant with the nuances of AI, while AI technologists might lack
953 knowledge of water treatment, supply, and distribution. This knowledge gap impedes the effective
954 deployment of ML in enhancing safety drinking water supply. Addressing this dichotomy is

955 beneficial and essential, necessitating educational and collaborative efforts to build a shared
956 understanding and to develop interdisciplinary skillsets.

957
958 Further complicating the matter is the absence of standardized toolkits tailored to safe drinking
959 water supply. Such standardization is vital for enabling consistent application across various
960 research and implementation efforts. Uniformity in tools and approaches would not only
961 streamline the processes but also bolster collaborative work, which is often fragmented across
962 regions and specializations.

963
964 Advancements in ML tools must cater to the unique challenges presented by safe drinking water
965 supply. Water quality is affected by numerous spatiotemporal factors, requiring ML solutions that
966 can adapt to and learn from these dynamic conditions. Thus, future studies should customize
967 existing ML frameworks or innovate new ones that can grapple with the complexities inherent to
968 safe drinking water supply.

969
970 Looking to the horizon, the broader vision involves leveraging ML to address the global drinking
971 water crisis. Issues such as water scarcity, the presence of emerging contaminants, and the
972 formation of DBPs present a global challenge. ML tools have been predominantly developed with
973 local or regional contexts, yet the drinking water crisis demands a global perspective. The ambition
974 to harness ML for these global challenges is critical to ensuring water security worldwide.

975
976 In pursuit of these goals, the integration of advanced ML models becomes a cornerstone in tackling
977 the multifaceted issues tied to drinking water safety. Future endeavors should prioritize the

978 promotion of open-access data sharing within and beyond the drinking water research community
979 (Zhong et al., 2021). The endeavors will enhance collaboration, drive transparency, and support
980 the reproducibility of scientific findings, which are the bedrock of robust research. Furthermore,
981 establishing a comprehensive comparative framework to evaluate different ML models will be
982 instrumental in identifying the optimal solutions for the challenges in drinking water research. By
983 embracing these strategies, we can aspire to not just bridge existing knowledge gaps but also
984 significantly elevate the role of ML in securing safe and more sustainable water supply.

985

986 9. Conclusions

987 Assessing and ensuring safe drinking water supply is a global challenge with conventional
988 approaches. ML as a novel tool is promising in monitoring and protecting drinking water quality,
989 especially in municipal engineered water systems. This review for the first time comprehensively
990 summarizes the applications of ML in assessing and ensuring safe drinking water supply with a
991 focus on water quality in engineered water systems. We compile the applications of ML from the
992 physical, chemical, microbiological, and temporal perspectives. From the physical perspective,
993 ML is useful in managing drinking water production and demand and monitoring drinking water
994 pipeline failures. From the chemical perspective, ML is promising in assessing and controlling
995 DBPs, monitoring and mitigating heavy metals, and tracking nitrification in drinking water. From
996 the microbiological perspective, ML can monitor and mitigate OPs, detect *Cryptosporidium* and
997 *Giardia*, assess biofilm development, assess AMR risks, and study microbial communities in
998 municipal water, especially in engineered water systems. In addition, ML is a useful tool in
999 assessing drinking water quality from the temporal perspective, especially in detecting accidental

1000 drinking water contamination. Taken together, ML is feasible in assessing and ensuring drinking
1001 water quality with a great potential to mitigate the global water crisis.

1002

1003 **Acknowledgment**

1004 This work has been supported by the Faculty Research Award from Arkansas State University
1005 (Jonesboro, AR, U.S.A.).

1006 **Declaration of interest**

1007 The authors declare that they have no known competing financial interests or personal
1008 relationships that could have appeared to influence the work reported in this paper.

1009

1010 References

- 1011 Abd Elnabi, M.K.; Elkaliny, N.E.; Elyazied, M.M.; Azab, S.H.; Elkhalfifa, S.A.; Elmasry, S.;
1012 Mouhamed, M.S.; Shalamesh, E.M.; Alhorieny, N.A.; Abd Elaty, A.E. **2023**. Toxicity of heavy
1013 metals and recent advances in their removal: a review. *Toxics* 11(7), 580.
- 1014 Abdi, J.; Mazloom, G. **2022**. Machine learning approaches for predicting arsenic adsorption from
1015 water using porous metal-organic frameworks. *Scientific Reports* 12, 16458.
- 1016 Abkar, L.; Moghaddam, H.S.; Fowler, S.J. **2024**. Microbial ecology of drinking water from source
1017 to tap. *Science of the Total Environment* 908, 168077.
- 1018 Adeyemo, F.E.; Singh, G.; Reddy, P.; Bux, F.; Stenstrom, T.A. **2019**. Efficiency of chlorine and
1019 UV in the inactivation of *Cryptosporidium* and *Giardia* in wastewater. *PLOS ONE* 14(5),
1020 e0216040.
- 1021 Ahmadpour, E.; Delpla, I.; Debia, M.; Simard, S.; Proulx, F.; Sérodes, J.-B.; Valois, I.; Tardif, R.;
1022 Haddad, S.; Rodriguez, M. **2023**. Full-scale multisampling and empirical modeling of DBPs
1023 in water and air of indoor pools. *Environmental Monitoring and Assessment* 195(9), 1128.
- 1024 Ahmed, A.N.; Othman, F.B.; Afan, H.A.; Ibrahim, R.K.; Fai, C.M.; Hossain, M.S.; Ehteram, M.;
1025 Elshafie, A. **2019**. Machine learning methods for better water quality prediction. *Journal of*
1026 *Hydrology* 578, 124084.
- 1027 Alexander, M.T.; Woodruff, P.; Mistry, J.H.; Buse, H.Y.; Lytle, D.A.; Pressman, J.G.; Wahman,
1028 D.G. **2024**. Evaluation of distribution system water quality during a free chlorine conversion
1029 [Manuscript submitted for publication].
- 1030 Alfredo, K. **2021**. The “Burn”: water quality and microbiological impacts related to limited free
1031 chlorine disinfection periods in a chloramine system. *Water Research* 197, 117044.
- 1032 Aliashrafi, A.; Zhang, Y.; Groenewegen, H.; Peleato, N.M. **2021**. A review of data-driven
1033 modelling in drinking water treatment. *Reviews in Environmental Science and Bio/Technology*
1034 20(4), 985-1009.
- 1035 Allen, J.M.; Plewa, M.J.; Wagner, E.D.; Wei, X.; Bokenkamp, K.; Hur, K.; Jia, A.; Liberatore,
1036 H.K.; Lee, C.-F.T.; Shirkhani, R. **2022**. Feel the burn: disinfection byproduct formation and
1037 cytotoxicity during chlorine burn events. *Environmental Science & Technology* 56(12), 8245–
1038 8254.
- 1039 Allesoe, R.L.; Lundgaard, A.T.; Hernandez Medina, R.; Aguayo-Orozco, A.; Johansen, J.; Nissen,
1040 J.N.; Brorsson, C.; Mazzoni, G.; Niu, L.; Biel, J.H.; Leal Rodriguez, C.; Brasas, V.; Webel, H.;
1041 Benros, M.E.; Pedersen, A.G.; Chmura, P.J.; Jacobsen, U.P.; Mari, A.; Koivula, R.; Mahajan,
1042 A.; Vinuela, A.; Tajés, J.F.; Sharma, S.; Haid, M.; Hong, M.G.; Musholt, P.B.; De Masi, F.;
1043 Vogt, J.; Pedersen, H.K.; Gudmundsdottir, V.; Jones, A.; Kennedy, G.; Bell, J.; Thomas, E.L.;
1044 Frost, G.; Thomsen, H.; Hansen, E.; Hansen, T.H.; Vestergaard, H.; Muilwijk, M.; Blom, M.T.;
1045 ‘t Hart, L.M.; Pattou, F.; Raverdy, V.; Brage, S.; Kokkola, T.; Heggie, A.; McEvoy, D.; Mourby,
1046 M.; Kaye, J.; Hattersley, A.; McDonald, T.; Ridderstrale, M.; Walker, M.; Forgie, I.; Giordano,
1047 G.N.; Pavo, I.; Ruetten, H.; Pedersen, O.; Hansen, T.; Dermitzakis, E.; Franks, P.W.; Schwenk,
1048 J.M.; Adamski, J.; McCarthy, M.I.; Pearson, E.; Banasik, K.; Rasmussen, S.; Brunak, S.;
1049 Consortium, I.D. **2023**. Discovery of drug-omics associations in type 2 diabetes with
1050 generative deep-learning models. *Nature Biotechnology* 41(3), 399-408.
- 1051 Almalioglu, Y.; Turan, M.; Trigoni, N.; Markham, A. **2022**. Deep learning-based robust positioning
1052 for all-weather autonomous driving. *Nature Machine Intelligence* 4(9), 749-760.
- 1053 Almheiri, Z.; Meguid, M.; Zayed, T. **2021**. Failure modeling of water distribution pipelines using
1054 meta-learning algorithms. *Water Research* 205, 117680.

- 1055 Ashbolt, N.J. **2015**. Microbial contamination of drinking water and human health from community
1056 water systems. *Current Environmental Health Reports* 2, 95-106.
- 1057 Asheri-Arnon, T.; Ezra, S.; Fishbain, B. **2018**. Contamination Detection of Water with Varying
1058 Routine Backgrounds by UV-Spectrophotometry. *Journal of Water Resources Planning and*
1059 *Management* 144(9), 04018056.
- 1060 Asheri Arnon, T.; Ezra, S.; Fishbain, B. **2019**. Water characterization and early contamination
1061 detection in highly varying stochastic background water, based on Machine Learning
1062 methodology for processing real-time UV-Spectrophotometry. *Water Research* 155, 333-342.
- 1063 Ates, N.; Civelekoglu, G.; Kaplan-Bekaroglu, S.S. **2022**. Management strategies for minimising
1064 DBPs formation in drinking water systems. In: Bahadir, M. and Haarstrick, A. (Eds.) *Water*
1065 *and Wastewater Management: Global Problems and Measures*. pp. 67-82. Springer
1066 International Publishing, Cham.
- 1067 Avni, N.; Fishbain, B.; Shamir, U. **2015**. Water consumption patterns as a basis for water demand
1068 modeling. *Water Resources Research* 51(10), 8165-8181.
- 1069 AWWA **2013**. *Nitrification Prevention and Control in Drinking Water*, 2nd Edition. American
1070 Water Works Association, Denver, Colorado, U.S.A.
- 1071 Bagriacik, A.; Davidson, R.A.; Hughes, M.W.; Bradley, B.A.; Cubrinovski, M. **2018**. Comparison
1072 of statistical and machine learning approaches to modeling earthquake damage to water
1073 pipelines. *Soil Dynamics and Earthquake Engineering* 112, 76-88.
- 1074 Bakker, M.; Vreeburg, J.H.G.; Van De Roer, M.; Rietveld, L.C. **2014**. Heuristic burst detection
1075 method using flow and pressure measurements. *Journal of Hydroinformatics* 16(5), 1194-1209.
- 1076 Balogun, A.-L.; Tella, A.; Baloo, L.; Adebisi, N. **2021**. A review of the inter-correlation of climate
1077 change, air pollution and urban sustainability using novel machine learning algorithms and
1078 spatial information science. *Urban Climate* 40, 100989.
- 1079 Benedict, K.M.; Reses, H.; Vigar, M.; Roth, D.M.; Roberts, V.A.; Mattioli, M.; Cooley, L.A.;
1080 Hilborn, E.D.; Wade, T.J.; Fullerton, K.E. **2017**. Surveillance for waterborne disease outbreaks
1081 associated with drinking water-United States, 2013-2014. *Morbidity and Mortality Weekly*
1082 *Report* 66(44), 1216.
- 1083 Benítez, J.S.; Rodríguez, C.M.; Casas, A.F. **2021**. Disinfection byproducts (DBPs) in drinking
1084 water supply systems: a systematic review. *Physics and Chemistry of the Earth, Parts A/B/C*
1085 123, 102987.
- 1086 Berglund, E.; Vizanko, B.; Kadinski, L.; Ostfeld, A. **2023**. Coupling machine learning and agent-
1087 based modeling to characterize contamination sources in water distribution systems for
1088 changing demand regimes. In: *World Environmental and Water Resources Congress 2023*. pp.
1089 10.1061/9780784484852.9780784484082.
- 1090 Berne, C.; Ellison, C.K.; Ducret, A.; Brun, Y.V. **2018**. Bacterial adhesion at the single-cell level.
1091 *Nature Reviews Microbiology* 16(10), 616-627.
- 1092 Bhagat, S.K.; Tung, T.M.; Yaseen, Z.M. **2020**. Development of artificial intelligence for modeling
1093 wastewater heavy metal removal: State of the art, application assessment and possible future
1094 research. *Journal of Cleaner Production* 250, 119473.
- 1095 Bi, K.; Xie, L.; Zhang, H.; Chen, X.; Gu, X.; Tian, Q. **2023**. Accurate medium-range global weather
1096 forecasting with 3D neural networks. *Nature* 619(7970), 533-538.
- 1097 Bond, T.; Goslan, E.H.; Parsons, S.A.; Jefferson, B. **2011**. Treatment of disinfection by-product
1098 precursors. *Environmental Technology* 32(1-2), 1-25.

- 1099 Bond, T.; Templeton, M.R.; Graham, N. **2012**. Precursors of nitrogenous disinfection by-products
1100 in drinking water—A critical review and analysis. *Journal of Hazardous Materials* 235-236,
1101 1-16.
- 1102 Bond, T.; Graham, N. **2017**. Predicting chloroform production from organic precursors. *Water*
1103 *Research* 124, 167-176.
- 1104 Bradley, A.P. **1997**. The use of the area under the ROC curve in the evaluation of machine learning
1105 algorithms. *Pattern Recognition* 30(7), 1145-1159.
- 1106 Breiman, L. **2001**. Random Forests. *Machine Learning* 45(1), 5-32.
- 1107 Brunello, A.; Civilini, M.; De Martin, S.; Felice, A.; Franchi, M.; Iacumin, L.; Saccomanno, N.;
1108 Vitacolonna, N. **2022**. Machine learning-assisted environmental surveillance of Legionella: A
1109 retrospective observational study in Friuli-Venezia Giulia region of Italy in the period 2002–
1110 2019. *Informatics in Medicine Unlocked* 28, 100803.
- 1111 Bruno, A.; Sandionigi, A.; Bernasconi, M.; Panio, A.; Labra, M.; Casiraghi, M. **2018**. Changes in
1112 the drinking water microbiome: Effects of water treatments along the flow of two drinking
1113 water treatment plants in a urbanized area, Milan (Italy). *Frontiers in Microbiology* 9, 2557.
- 1114 CDC **2021a**. Parasites - *Cryptosporidium* (also known as "Crypto").
- 1115 CDC **2021b**. Parasites - *Giardia*.
- 1116 Cha, D.; Park, S.; Kim, M.S.; Kim, T.; Hong, S.W.; Cho, K.H.; Lee, C. **2021**. Prediction of oxidant
1117 exposures and micropollutant abatement during ozonation using a machine learning method.
1118 *Environmental Science & Technology* 55(1), 709-718.
- 1119 Chang, Q.; Wang, W.; Regev-Yochay, G.; Lipsitch, M.; Hanage, W.P. **2015**. Antibiotics in
1120 agriculture and the risk to human health: how worried should we be? *Evolutionary Applications*
1121 8(3), 240-247.
- 1122 Chen, H.; Bai, X.; Li, Y.; Jing, L.; Chen, R.; Teng, Y. **2019**. Source identification of antibiotic
1123 resistance genes in a peri-urban river using novel crAssphage marker genes and metagenomic
1124 signatures. *Water Research* 167, 115098.
- 1125 Chowdhury, S.; Mazumder, M.A.J.; Al-Attas, O.; Husain, T. **2016**. Heavy metals in drinking water:
1126 occurrences, implications, and future needs in developing countries. *Science of the Total*
1127 *Environment* 569-570, 476-488.
- 1128 Cordero, J.A.; He, K.; Janya, K.; Echigo, S.; Itoh, S. **2021**. Predicting formation of haloacetic acids
1129 by chlorination of organic compounds using machine-learning-assisted quantitative structure-
1130 activity relationships. *Journal of Hazardous Materials* 408, 124466.
- 1131 Craun, G.F.; Brunkard, J.M.; Yoder, J.S.; Roberts, V.A.; Carpenter, J.; Wade, T.; Calderon, R.L.;
1132 Roberts, J.M.; Beach, M.J.; Roy, S.L. **2010**. Causes of outbreaks associated with drinking water
1133 in the United States from 1971 to 2006. *Clinical Microbiology Reviews* 23(3), 507-528.
- 1134 Daniel, I.; Abhijith, G.R.; Kadinski, L.; Ostfeld, A.; Cominola, A. **2023**. A machine learning-based
1135 surrogate model for coupled hydraulic and water quality simulation in water distribution
1136 networks. In: *World Environmental and Water Resources Congress 2023*. pp.
1137 10.1061/9780784484852.9780784484077.
- 1138 De Jesus, K.L.M.; Senoro, D.B.; Dela Cruz, J.C.; Chan, E.B. **2021**. A hybrid neural network-
1139 particle swarm optimization informed spatial interpolation technique for groundwater quality
1140 mapping in a small island province of the Philippines. *Toxics* 9(11), 273.
- 1141 Delpla, I.; Jung, A.-V.; Baures, E.; Clement, M.; Thomas, O. **2009**. Impacts of climate change on
1142 surface water quality in relation to drinking water production. *Environment International* 35(8),
1143 1225-1233.

- 1144 Deng, Y.; Zhou, X.; Shen, J.; Xiao, G.; Hong, H.; Lin, H.; Wu, F.; Liao, B.-Q. **2021**. New methods
1145 based on back propagation (BP) and radial basis function (RBF) artificial neural networks
1146 (ANNs) for predicting the occurrence of haloketones in tap water. *Science of the Total*
1147 *Environment* 772, 145534.
- 1148 Dogo, E.M.; Nwulu, N.I.; Twala, B.; Aigbavboa, C. **2019**. A survey of machine learning methods
1149 applied to anomaly detection on drinking-water quality data. *Urban Water Journal* 16(3), 235-
1150 248.
- 1151 Donkor, E.A.; Mazzuchi, T.A.; Soyer, R.; Roberson, J.A. **2014**. Urban water demand forecasting:
1152 Review of methods and models. *Journal of Water Resources Planning and Management* 140(2),
1153 146-159.
- 1154 Donohue, M.; King, D.; Pfaller, S.; Mistry, J. **2019**. The sporadic nature of *Legionella pneumophila*,
1155 *Legionella pneumophila* Sg1 and *Mycobacterium avium* occurrence within residences and
1156 office buildings across 36 states in the United States. *Journal of Applied Microbiology* 126(5),
1157 1568-1579.
- 1158 Ewuzie, U.; Bolade, O.P.; Egbedina, A.O. **2022**. Application of deep learning and machine learning
1159 methods in water quality modeling and prediction: a Review. In: Marques, G. and Ighalo, J.O.
1160 (Eds.) *Current Trends and Advances in Computer-Aided Intelligent Environmental Data*
1161 *Engineering*. pp. 185-218. Elsevier Inc., San Diego, California, U.S.A.
- 1162 Fan, X.; Zhang, X.; Yu, A.; Speitel, M.; Yu, X. **2023**. Assessment of the impacts of climate change
1163 on water supply system pipe failures. *Scientific Reports* 13(1), 7349.
- 1164 Fasel, U.; Kutz, J.N.; Brunton, B.W.; Brunton, S.L. **2022**. Ensemble-SINDy: Robust sparse model
1165 discovery in the low-data, high-noise limit, with active learning and control. *Proceedings of*
1166 *the Royal Society A: Mathematical, Physical and Engineering Sciences* 478(2260), 20210904.
- 1167 Favere, J.; Barbosa, R.G.; Sleutels, T.; Verstraete, W.; De Gussemme, B.; Boon, N. **2021**.
1168 Safeguarding the microbial water quality from source to tap. *NPJ Clean Water* 4, 28.
- 1169 Feng, S.; Sun, H.; Yan, X.; Zhu, H.; Zou, Z.; Shen, S.; Liu, H.X. **2023**. Dense reinforcement
1170 learning for safety validation of autonomous vehicles. *Nature* 615(7953), 620-627.
- 1171 Figueras, M.J.; Borrego, J.J. **2010**. New perspectives in monitoring drinking water microbial
1172 quality. *International Journal of Environmental Research and Public Health* 7(12), 4179-4202.
- 1173 Flemming, H.-C.; Wingender, J. **2010**. The biofilm matrix. *Nature Reviews Microbiology* 8(9),
1174 623-633.
- 1175 Flemming, H.-C.; van Hullebusch, E.D.; Neu, T.R.; Nielsen, P.H.; Seviour, T.; Stoodley, P.;
1176 Wingender, J.; Wuertz, S. **2023**. The biofilm matrix: multitasking in a shared space. *Nature*
1177 *Reviews Microbiology* 21(2), 70-86.
- 1178 Friedman, J.H. **2001**. Greedy function approximation: a gradient boosting machine. *The Annals of*
1179 *Statistics* 29(5), 1189-1232.
- 1180 Fu, Z.; Xi, S. **2020**. The effects of heavy metals on human metabolism. *Toxicology mechanisms*
1181 *and methods* 30(3), 167-176.
- 1182 Ghiassi, M.; Fa'al, F.; Abrishamchi, A. **2017**. Large metropolitan water demand forecasting using
1183 DAN2, FTDNN, and KNN models: A case study of the city of Tehran, Iran. *Urban Water*
1184 *Journal* 14(6), 655-659.
- 1185 Ghobadi, F.; Kang, D. **2023**. Application of machine learning in water resources management: a
1186 systematic literature review. *Water* 15(4), 620.
- 1187 Gogoi, A.; Mazumder, P.; Tyagi, V.K.; Chaminda, G.T.; An, A.K.; Kumar, M. **2018**. Occurrence
1188 and fate of emerging contaminants in water environment: a review. *Groundwater for*
1189 *Sustainable Development* 6, 169-180.

- 1190 Gomez-Alvarez, V.; Revetta, R.P. **2020**. Monitoring of nitrification in chloraminated drinking
1191 water distribution systems with microbiome bioindicators using supervised machine learning.
1192 *Frontiers in Microbiology* 11, 2254.
- 1193 Gong, J.; Guo, X.; Yan, X.; Hu, C. **2023**. Review of urban drinking water contamination source
1194 identification methods. *Energies* 16(2), 705.
- 1195 Gonzalez, S.; Lopez-Roldan, R.; Cortina, J.-L. **2013**. Presence of metals in drinking water
1196 distribution networks due to pipe material leaching: a review. *Toxicological & Environmental*
1197 *Chemistry* 95(6), 870-889.
- 1198 Goodfellow, I.; Bengio, Y.; Courville, A. **2016**. *Deep Learning*. MIT Press.
- 1199 Greenberg, J.; Price, B.; Ware, A. **2010**. Alternative estimate of source distribution in microbial
1200 source tracking using posterior probabilities. *Water Research* 44(8), 2629-2637.
- 1201 Grey, D.; Sadoff, C.W. **2007**. Sink or swim? Water security for growth and development. *Water*
1202 *Policy* 9(6), 545-571.
- 1203 Guo, G.; Liu, S.; Wu, Y.; Li, J.; Zhou, R.; Zhu, X. **2018**. Short-term water demand forecast based
1204 on deep learning method. *Journal of Water Resources Planning and Management* 144(12),
1205 04018076.
- 1206 Gupta, S.; Basant, N. **2016**. Modeling the reactivity of ozone and sulphate radicals towards organic
1207 chemicals in water using machine learning approaches. *RSC advances* 6(110), 108448-108457.
- 1208 Hamidian, A.H.; Esfandeh, S.; Zhang, Y.; Yang, M. **2019**. Simulation and optimization of
1209 nanomaterials application for heavy metal removal from aqueous solutions. *Inorganic and*
1210 *Nano-Metal Chemistry* 49(7), 217-230.
- 1211 Hastie, T.; Tibshirani, R.; Friedman, J.H. **2009**. *The Elements of Statistical Learning: Data Mining,*
1212 *Inference, and Prediction*, 2nd Edition. Springer, New York.
- 1213 He, L.; Bai, L.; Dionysiou, D.D.; Wei, Z.; Spinney, R.; Chu, C.; Lin, Z.; Xiao, R. **2021**.
1214 Applications of computational chemistry, artificial intelligence, and machine learning in
1215 aquatic chemistry research. *Chemical Engineering Journal* 426, 131810.
- 1216 Helte, E.; Säve-Söderbergh, M.; Larsson, S.C.; Martling, A.; Åkesson, A. **2023**. Disinfection by-
1217 products in drinking water and risk of colorectal cancer: a population-based cohort study.
1218 *Journal of the National Cancer Institute* 115(12), 1597-1604.
- 1219 Henrique Alves Ribeiro, V.; Reynoso-Meza, G. **2023**. Multi-criteria decision-making techniques
1220 for the selection of Pareto-optimal machine learning models in a drinking-water quality
1221 monitoring problem. *International Journal of Information Technology & Decision Making*
1222 23(1), 447-474.
- 1223 Hong, H.; Zhang, Z.; Guo, A.; Shen, L.; Sun, H.; Liang, Y.; Wu, F.; Lin, H. **2020**. Radial basis
1224 function artificial neural network (RBF ANN) as well as the hybrid method of RBF ANN and
1225 grey relational analysis able to well predict trihalomethanes levels in tap water. *Journal of*
1226 *Hydrology* 591, 125574.
- 1227 Hossain, S.; Cook, D.; Chow, C.W.K.; Hewa, G.A. **2021**. Development of an optical method to
1228 monitor nitrification in drinking water. *Sensors* 21(22), 7525.
- 1229 Hu, G.; Mian, H.R.; Mohammadiun, S.; Rodriguez, M.J.; Hewage, K.; Sadiq, R. **2023**. Appraisal
1230 of machine learning techniques for predicting emerging disinfection byproducts in small water
1231 distribution networks. *Journal of Hazardous Materials* 446, 130633.
- 1232 Huang, R.; Ma, C.; Ma, J.; Huangfu, X.; He, Q. **2021**. Machine learning in natural and engineered
1233 water systems. *Water Research* 205, 117666.

- 1234 Huang, Y.; Li, T.; Zheng, S.; Fan, L.; Su, L.; Zhao, Y.; Xie, H.B.; Li, C. **2020**. QSAR modeling for
1235 the ozonation of diverse organic compounds in water. *Science of the Total Environment* 715,
1236 136816.
- 1237 Hutchings, M.I.; Truman, A.W.; Wilkinson, B. **2019**. Antibiotics: past, present and future. *Current*
1238 *Opinion in Microbiology* 51, 72-80.
- 1239 Isaac, T.S.; Sherchan, S.P. **2020**. Molecular detection of opportunistic premise plumbing pathogens
1240 in rural Louisiana's drinking water distribution system. *Environmental Research* 181, 108847.
- 1241 Jefri, U.H.N.M.; Khan, A.; Lim, Y.C.; Lee, K.S.; Liew, K.B.; Kassab, Y.W.; Choo, C.-Y.; Al-Worafi,
1242 Y.M.; Ming, L.C.; Kalusalingam, A. **2022**. A systematic review on chlorine dioxide as a
1243 disinfectant. *Journal of Medicine and Life* 15(3), 313.
- 1244 Jelli, E.; Ohmura, T.; Netter, N.; Abt, M.; Jiménez-Siebert, E.; Neuhaus, K.; Rode, D.K.H.; Nadell,
1245 C.D.; Drescher, K. **2023**. Single-cell segmentation in bacterial biofilms with an optimized deep
1246 learning method enables tracking of cell lineages and measurements of growth rates. *Molecular*
1247 *Microbiology* 119(6), 659-676.
- 1248 Jia, X.; O'Connor, D.; Shi, Z.; Hou, D. **2021**. VIRS based detection in combination with machine
1249 learning for mapping soil pollution. *Environ Pollut* 268(Pt A), 115845.
- 1250 Jiang, S.; Hu, J.; Wood, K.L.; Luo, J. **2021**. Data-driven design-by-analogy: state-of-the-art and
1251 future directions. *Journal of Mechanical Design* 144(2), 020801.
- 1252 Jing, Z.; Lu, Z.; Zhao, Z.; Cao, W.; Wang, W.; Ke, Y.; Wang, X.; Sun, W. **2023**. Molecular
1253 ecological networks reveal the spatial-temporal variation of microbial communities in drinking
1254 water distribution systems. *Journal of Environmental Sciences* 124, 176-186.
- 1255 Joseph, L.; Jun, B.-M.; Flora, J.R.V.; Park, C.M.; Yoon, Y. **2019**. Removal of heavy metals from
1256 water sources in the developing world using low-cost materials: A review. *Chemosphere* 229,
1257 142-159.
- 1258 Karygianni, L.; Ren, Z.; Koo, H.; Thurnheer, T. **2020**. Biofilm matrixome: extracellular
1259 components in structured microbial communities. *Trends in Microbiology* 28(8), 668-681.
- 1260 Kazemi, E.; Kyritsakas, G.; Husband, S.; Flavell, K.; Speight, V.; Boxall, J. **2023**. Predicting iron
1261 exceedance risk in drinking water distribution systems using machine learning. *IOP*
1262 *Conference Series: Earth and Environmental Science* 1136(1), 012047.
- 1263 Khattak, A.; Bukhsh, R.; Aslam, S.; Yafoz, A.; Alghushairy, O.; Alsini, R. **2022**. A hybrid deep
1264 learning-based model for detection of electricity losses using big data in power systems.
1265 *Sustainability* 14(20), 13627.
- 1266 Kirstein, I.V.; Hensel, F.; Gomiero, A.; Iordachescu, L.; Vianello, A.; Wittgren, H.B.; Vollertsen, J.
1267 **2021**. Drinking plastics?—Quantification and qualification of microplastics in drinking water
1268 distribution systems by μ FTIR and Py-GCMS. *Water Research* 188, 116519.
- 1269 Knights, D.; Kuczynski, J.; Charlson, E.S.; Zaneveld, J.; Mozer, M.C.; Collman, R.G.; Bushman,
1270 F.D.; Knight, R.; Kelley, S.T. **2011**. Bayesian community-wide culture-independent microbial
1271 source tracking. *Nature Methods* 8(9), 761-763.
- 1272 Krasner, S.W.; Weinberg, H.S.; Richardson, S.D.; Pastor, S.J.; Chinn, R.; Scilimenti, M.J.; Onstad,
1273 G.D.; Thruston, A.D. **2006**. Occurrence of a New Generation of Disinfection Byproducts.
1274 *Environmental Science & Technology* 40(23), 7175-7185.
- 1275 Krasner, S.W.; Mitch, W.A.; McCurry, D.L.; Hanigan, D.; Westerhoff, P. **2013**. Formation,
1276 precursors, control, and occurrence of nitrosamines in drinking water: a review. *Water*
1277 *Research* 47(13), 4433-4450.

- 1278 Lambert, A.; Brown, T.G.; Takizawa, M.; Weimer, D. **1999**. A review of performance indicators
1279 for real losses from water supply systems. *Journal of Water Supply: Research and*
1280 *Technology—AQUA* 48(6), 227-237.
- 1281 LeCun, Y.; Boser, B.; Denker, J.S.; Henderson, D.; Howard, R.E.; Hubbard, W.; Jackel, L.D. **1989**.
1282 Backpropagation Applied to Handwritten Zip Code Recognition. *Neural Computation* 1(4),
1283 541-551.
- 1284 Lecun, Y.; Bottou, L.; Bengio, Y.; Haffner, P. **1998**. Gradient-based learning applied to document
1285 recognition. *Proceedings of the IEEE* 86(11), 2278-2324.
- 1286 Lee, D.; Gibson, J.M.; Brown, J.; Habtewold, J.; Murphy, H.M. **2023**. Burden of disease from
1287 contaminated drinking water in countries with high access to safely managed water: a
1288 systematic review. *Water Research* 242, 120244.
- 1289 Lee, E.J.; Schwab, K.J. **2005**. Deficiencies in drinking water distribution systems in developing
1290 countries. *Journal of Water and Health* 3(2), 109-127.
- 1291 Lee, J.; Kim, E.-S.; Roh, B.-S.; Eom, S.-W.; Zoh, K.-D. **2013**. Occurrence of disinfection by-
1292 products in tap water distribution systems and their associated health risk. *Environmental*
1293 *Monitoring and Assessment* 185, 7675-7691.
- 1294 Legube, B.; Parinet, B.; Gelinet, K.; Berne, F.; Croue, J.P. **2004**. Modeling of bromate formation
1295 by ozonation of surface waters in drinking water treatment. *Water Research* 38(8), 2185-2195.
- 1296 Leitão, J.; Simões, N.; Sá Marques, J.A.; Gil, P.; Ribeiro, B.; Cardoso, A. **2019**. Detecting urban
1297 water consumption patterns: a time-series clustering approach. *Water Supply* 19(8), 2323-2329.
- 1298 Ley, C.; Martin, R.K.; Pareek, A.; Groll, A.; Seil, R.; Tischer, T. **2022**. Machine learning and
1299 conventional statistics: making sense of the differences. *Knee Surg Sports Traumatol Arthrosc*
1300 30(3), 753-757.
- 1301 Li, J.; Heap, A.D.; Potter, A.; Daniell, J.J. **2011**. Application of machine learning methods to spatial
1302 interpolation of environmental variables. *Environmental Modelling & Software* 26(12), 1647-
1303 1659.
- 1304 Li, L.; Rong, S.; Wang, R.; Yu, S. **2021**. Recent advances in artificial intelligence and machine
1305 learning for nonlinear relationship analysis and process control in drinking water treatment: A
1306 review. *Chemical Engineering Journal* 405, 126673.
- 1307 Li, P.; Wu, J. **2019**. Drinking water quality and public health. *Exposure and Health* 11(2), 73-79.
- 1308 Li, Q.; Yu, S.; Li, L.; Liu, G.; Gu, Z.; Liu, M.; Liu, Z.; Ye, Y.; Xia, Q.; Ren, L. **2017**. Microbial
1309 communities shaped by treatment processes in a drinking water treatment plant and their
1310 contribution and threat to drinking water safety. *Frontiers in Microbiology* 8, 2465.
- 1311 Li, R.A.; McDonald, J.A.; Sathasivan, A.; Khan, S.J. **2019**. Disinfectant residual stability leading
1312 to disinfectant decay and by-product formation in drinking water distribution systems: a
1313 systematic review. *Water Research* 153, 335-348.
- 1314 Li, Z.; Zhang, C.; Liu, H.; Zhang, C.; Zhao, M.; Gong, Q.; Fu, G. **2022**. Developing stacking
1315 ensemble models for multivariate contamination detection in water distribution systems.
1316 *Science of the Total Environment* 828, 154284.
- 1317 Li, Z.; Liu, H.; Zhang, C.; Fu, G. **2023**. Generative adversarial networks for detecting
1318 contamination events in water distribution systems using multi-parameter, multi-site water
1319 quality monitoring. *Environmental Science and Ecotechnology* 14, 100231.
- 1320 Ligda, P.; Claerebout, E.; Kostopoulou, D.; Zdragas, A.; Casaert, S.; Robertson, L.J.; Sotiraki, S.
1321 **2020**. *Cryptosporidium* and *Giardia* in surface water and drinking water: Animal sources and
1322 towards the use of a machine-learning approach as a tool for predicting contamination.
1323 *Environmental Pollution* 264, 114766.

- 1324 Lin, H.; Dai, Q.; Zheng, L.; Hong, H.; Deng, W.; Wu, F. **2020**. Radial basis function artificial
1325 neural network able to accurately predict disinfection by-product levels in tap water: Taking
1326 haloacetic acids as a case study. *Chemosphere* 248, 125999.
- 1327 Lin, H.H.-H.; Lin, A.Y.-C. **2024**. Peracetic acid as an alternative disinfectant for micropollutants
1328 degradation and disinfection byproducts control in outdoor swimming pools. *Journal of*
1329 *Hazardous Materials* 464, 132988.
- 1330 Ling, F.; Whitaker, R.; LeChevallier, M.W.; Liu, W.T. **2018**. Drinking water microbiome assembly
1331 induced by water stagnation. *The ISME Journal* 12(6), 1520-1531.
- 1332 Liu, G.; Zhang, Y.; van der Mark, E.; Magic-Knezev, A.; Pinto, A.; van den Bogert, B.; Liu, W.;
1333 van der Meer, W.; Medema, G. **2018**. Assessing the origin of bacteria in tap water and
1334 distribution system in an unchlorinated drinking water system by SourceTracker using
1335 microbial community fingerprints. *Water Research* 138, 86-96.
- 1336 Liu, S.; Gunawan, C.; Barraud, N.; Rice, S.; Harry, E.; Amal, R. **2016**. Understanding, monitoring
1337 and controlling biofilm growth in drinking water distribution systems. *Environmental Science*
1338 *& Technology* 50(17), 8954-8976.
- 1339 Lowe, M.; Qin, R.; Mao, X. **2022**. A review on machine learning, artificial intelligence, and smart
1340 technology in water treatment and monitoring. *Water* 14(9), 1384.
- 1341 Lozupone, C.; Lladser, M.E.; Knights, D.; Stombaugh, J.; Knight, R. **2011**. UniFrac: an effective
1342 distance metric for microbial community comparison. *The ISME Journal* 5(2), 169-172.
- 1343 Luo, S.; Nguyen, K.T.; Nguyen, B.T.T.; Feng, S.; Shi, Y.; Elsayed, A.; Zhang, Y.; Zhou, X.; Wen,
1344 B.; Chierchia, G.; Talbot, H.; Bourouina, T.; Jiang, X.; Liu, A.Q. **2021**. Deep learning-enabled
1345 imaging flow cytometry for high-speed *Cryptosporidium* and *Giardia* detection. *Cytometry*
1346 *Part A* 99(11), 1123-1133.
- 1347 Lytle, D.A.; Pfaller, S.; Muhlen, C.; Struewing, I.; Triantafyllidou, S.; White, C.; Hayes, S.; King,
1348 D.; Lu, J. **2021**. A comprehensive evaluation of monochloramine disinfection on water quality,
1349 *Legionella* and other important microorganisms in a hospital. *Water Research* 189, 116656.
- 1350 Mahajna, A.; Dinkla, I.J.; Euverink, G.J.W.; Keesman, K.J.; Jayawardhana, B. **2022**. Clean and
1351 safe drinking water systems via metagenomics data and artificial intelligence: state-of-the-art
1352 and future perspective. *Frontiers in Microbiology* 13, 832452.
- 1353 Manasfi, T. **2021**. Ozonation in drinking water treatment: an overview of general and practical
1354 aspects, mechanisms, kinetics, and byproduct formation. *Comprehensive Analytical Chemistry*
1355 92, 85-116.
- 1356 Mao, Y.; Wang, X.; Yang, H.; Wang, H.; Xie, Y.F. **2014**. Effects of ozonation on disinfection
1357 byproduct formation and speciation during subsequent chlorination. *Chemosphere* 117, 515-
1358 520.
- 1359 Mays, L.W. **2000**. *Water Distribution System Handbook*, 1st Edition. McGraw-Hill Education,
1360 New York.
- 1361 Mazhar, M.A.; Khan, N.A.; Ahmed, S.; Khan, A.H.; Hussain, A.; Rahisuddin; Changani, F.;
1362 Yousefi, M.; Ahmadi, S.; Vambol, V. **2020**. Chlorination disinfection by-products in municipal
1363 drinking water – A review. *Journal of Cleaner Production* 273, 123159.
- 1364 Mounce, S.R.; Mounce, R.B.; Boxall, J.B. **2010**. Novelty detection for time series data analysis in
1365 water distribution systems using support vector machines. *Journal of Hydroinformatics* 13(4),
1366 672-686.
- 1367 Muharemi, F.; Logofătu, D.; Leon, F. **2019**. Machine learning approaches for anomaly detection
1368 of water quality on a real-world data set. *Journal of Information and Telecommunication* 3(3),
1369 294-307.

- 1370 Mukhopadhyay, A.; Duttagupta, S.; Mukherjee, A. **2022**. Emerging organic contaminants in global
1371 community drinking water sources and supply: A review of occurrence, processes and
1372 remediation. *Journal of Environmental Chemical Engineering* 10(3), 107560.
- 1373 Naloufi, M.; Lucas, F.S.; Souihi, S.; Servais, P.; Janne, A.; Wanderley Matos De Abreu, T. **2021**.
1374 Evaluating the performance of machine learning approaches to predict the microbial quality of
1375 surface waters and to optimize the sampling effort. *Water* 13(18), 2457.
- 1376 Narita, K.; Matsui, Y.; Matsushita, T.; Shirasaki, N. **2023**. Screening priority pesticides for drinking
1377 water quality regulation and monitoring by machine learning: Analysis of factors affecting
1378 detectability. *Journal of environmental management* 326, 116738.
- 1379 Oh, S.; Hossen, I.; Luglio, J.; Justin, G.; Richie, J.E.; Medeiros, H.; Lee, C.H. **2021**. *On-site/in situ*
1380 continuous detecting ppb-level metal ions in drinking water using block loop-gap resonators
1381 and machine learning. *IEEE Transactions on Instrumentation and Measurement* 70, 9513909.
- 1382 Olikier, N.; Ostfeld, A. **2014**. Minimum volume ellipsoid classification model for contamination
1383 event detection in water distribution systems. *Environmental Modelling & Software* 57, 1-12.
- 1384 Ortiz-Lopez, C.; Bouchard, C.; Rodriguez, M. **2022**. Machine learning models with potential
1385 application to predict source water quality for treatment purposes: a critical review.
1386 *Environmental Technology Reviews* 11(1), 118-147.
- 1387 Pan, R.; Zhang, T.-Y.; Zheng, Z.-X.; Ai, J.; Ye, T.; Zhao, H.-X.; Hu, C.-Y.; Tang, Y.-L.; Fan, J.-J.;
1388 Geng, B.; Xu, B. **2023**. Insight into mixed chlorine/chloramines conversion and associated
1389 water quality variability in drinking water distribution systems. *Science of the Total*
1390 *Environment* 880, 163297.
- 1391 Pandian, A.M.K.; Rajamehala, M.; Singh, M.V.P.; Sarojini, G.; Rajamohan, N. **2022**. Potential
1392 risks and approaches to reduce the toxicity of disinfection by-product—A review. *Science of the*
1393 *Total Environment* 822, 153323.
- 1394 Pang, G.; Cao, L.; Chen, L.; Lian, D.; Liu, H. **2018**. Sparse modeling-based sequential ensemble
1395 learning for effective outlier detection in high-dimensional numeric data. *Thirty-Second AAAI*
1396 *Conference on Artificial Intelligence* 32(1), 10.1609/aaai.v1632i1601.11692.
- 1397 Park, J.; Park, J.-H.; Choi, J.-S.; Joo, J.C.; Park, K.; Yoon, H.C.; Park, C.Y.; Lee, W.H.; Heo, T.-Y.
1398 **2020**. Ensemble model development for the prediction of a disaster index in water treatment
1399 systems. *Water* 12(11), 3195.
- 1400 Park, M.; Anumol, T.; Snyder, S.A. **2015**. Modeling approaches to predict removal of trace organic
1401 compounds by ozone oxidation in potable reuse applications. *Environmental Science: Water*
1402 *Research & Technology* 1(5), 699-708.
- 1403 Peleato, N.M.; Legge, R.L.; Andrews, R.C. **2018**. Neural networks for dimensionality reduction of
1404 fluorescence spectra and prediction of drinking water disinfection by-products. *Water Research*
1405 136, 84-94.
- 1406 Peleato, N.M. **2022**. Application of convolutional neural networks for prediction of disinfection
1407 by-products. *Scientific Reports* 12(1), 612.
- 1408 Pifer, A.D.; Fairey, J.L. **2012**. Improving on SUVA₂₅₄ using fluorescence-PARAFAC analysis and
1409 asymmetric flow-field flow fractionation for assessing disinfection byproduct formation and
1410 control. *Water Research* 46(9), 2927-2936.
- 1411 Pinto, A.J.; Schroeder, J.; Lunn, M.; Sloan, W.; Raskin, L. **2014**. Spatial-temporal survey and
1412 occupancy-abundance modeling to predict bacterial community dynamics in the drinking
1413 water microbiome. *MBio* 5(3), e01135-14.
- 1414 Podgorski, J.; Berg, M. **2022**. Global analysis and prediction of fluoride in groundwater. *Nature*
1415 *Communications* 13, 4232.

- 1416 Prescott, J.F. **2017**. History and current use of antimicrobial drugs in veterinary medicine.
1417 *Microbiology Spectrum* 5(6), ARBA-0002-2017.
- 1418 Proctor, C.R.; Lee, J.; Yu, D.; Shah, A.D.; Whelton, A.J. **2020**. Wildfire caused widespread
1419 drinking water distribution network contamination. *AWWA Water Science* 2(4), e1183.
- 1420 Ramos-Martínez, E.; Herrera, M.; Izquierdo, J.; Pérez-García, R. **2014**. Ensemble of naïve
1421 Bayesian approaches for the study of biofilm development in drinking water distribution
1422 systems. *International Journal of Computer Mathematics* 91(1), 135-146.
- 1423 Ramos-Martínez, E.; Herrera, M.; Izquierdo, J.; Pérez-García, R. **2016**. A multi-disciplinary
1424 procedure to ascertain biofilm formation in drinking water pipes. *International Congress on*
1425 *Environmental Modelling and Software* 3, 619-626.
- 1426 Redondo-Hasselerharm, P.E.; Cserbik, D.; Flores, C.; Farré, M.J.; Sanchís, J.; Alcolea, J.A.; Planas,
1427 C.; Caixach, J.; Villanueva, C.M. **2022**. Insights to estimate exposure to regulated and non-
1428 regulated disinfection by-products in drinking water. *Journal of exposure science &*
1429 *environmental epidemiology* 34, 23-33.
- 1430 Reis, A.L.; Lopes, M.A.; Andrade-Campos, A.; Antunes, C.H. **2023**. A review of operational
1431 control strategies in water supply systems for energy and cost efficiency. *Renewable and*
1432 *Sustainable Energy Reviews* 175, 113140.
- 1433 Renwick, D.V.; Heinrich, A.; Weisman, R.; Arvanaghi, H.; Rotert, K. **2019**. Potential public health
1434 impacts of deteriorating distribution system infrastructure. *Journal AWWA* 111(2), 42-53.
- 1435 Richards, C.E.; Tzachor, A.; Avin, S.; Fenner, R. **2023**. Rewards, risks and responsible deployment
1436 of artificial intelligence in water systems. *Nature Water* 1(5), 422-432.
- 1437 Richardson, S.D.; Thruston, A.D.; Caughran, T.V.; Chen, P.H.; Collette, T.W.; Floyd, T.L.; Schenck,
1438 K.M.; Lykins, B.W.; Sun, G.-r.; Majetich, G. **1999**. Identification of new ozone disinfection
1439 byproducts in drinking water. *Environmental Science & Technology* 33(19), 3368-3377.
- 1440 Ritter, K.J.; Carruthers, E.; Carson, C.A.; Ellender, R.D.; Harwood, V.J.; Kingsley, K.; Nakatsu,
1441 C.; Sadowsky, M.; Shear, B.; West, B.; Whitlock, J.E.; Wiggins, B.A.; Wilbur, J.D. **2003**.
1442 Assessment of statistical methods used in library-based approaches to microbial source
1443 tracking. *Journal of Water and Health* 1(4), 209-223.
- 1444 Roca, I.; Akova, M.; Baquero, F.; Carlet, J.; Cavaleri, M.; Coenen, S.; Cohen, J.; Findlay, D.;
1445 Gyssens, I.; Heuer, O.E.; Kahlmeter, G.; Kruse, H.; Laxminarayan, R.; Liébana, E.; López-
1446 Cerero, L.; MacGowan, A.; Martins, M.; Rodríguez-Baño, J.; Rolain, J.M.; Segovia, C.;
1447 Sigauque, B.; Tacconelli, E.; Wellington, E.; Vila, J. **2015**. The global threat of antimicrobial
1448 resistance: science for intervention. *New Microbes New Infect* 6, 22-29.
- 1449 Rodríguez-Perez, J.; Leigh, C.; Liquet, B.; Kermorvant, C.; Peterson, E.; Sous, D.; Mengersen, K.
1450 **2020**. Detecting technical anomalies in high-frequency water-quality data using artificial
1451 neural networks. *Environmental Science & Technology* 54(21), 13719-13730.
- 1452 Rodríguez, M.J.; Sérodes, J.B.; Levallois, P. **2004**. Behavior of trihalomethanes and haloacetic
1453 acids in a drinking water distribution system. *Water Research* 38(20), 4367-4382.
- 1454 Saboe, D.; Ghasemi, H.; Gao, M.M.; Samardzic, M.; Hristovski, K.D.; Boscovic, D.; Burge, S.R.;
1455 Burge, R.G.; Hoffman, D.A. **2021**. Real-time monitoring and prediction of water quality
1456 parameters and algae concentrations using microbial potentiometric sensor signals and
1457 machine learning tools. *Science of the Total Environment* 764, 142876.
- 1458 Senoro, D.B.; de Jesus, K.L.M.; Nolos, R.C.; Lamac, M.R.L.; Deseo, K.M.; Tabelin, C.B. **2022**.
1459 In situ measurements of domestic water quality and health risks by elevated concentration of
1460 heavy metals and metalloids using Monte Carlo and MLGI methods. *Toxics* 10(7), 342.

1461 Shao, B.; Shen, L.; Liu, Z.; Tang, L.; Tan, X.; Wang, D.; Zeng, W.; Wu, T.; Pan, Y.; Zhang, X. **2023**.
1462 Disinfection byproducts formation from emerging organic micropollutants during chlorine-
1463 based disinfection processes. *Chemical Engineering Journal* 455, 140476.
1464 Shi, Y.; Babatunde, A.; Bockelmann-Evans, B.; Li, Q.; Zhang, L. **2020**. On-going nitrification in
1465 chloraminated drinking water distribution system (DWDS) is conditioned by hydraulics and
1466 disinfection strategies. *Journal of Environmental Sciences* 96, 151-162.
1467 Shi, Y.; Wang, J.; Wang, Q.; Jia, Q.; Yan, F.; Luo, Z.-H.; Zhou, Y.-N. **2022**. Supervised machine
1468 learning algorithms for predicting rate constants of ozone reaction with micropollutants.
1469 *Industrial & Engineering Chemistry Research* 61(24), 8359-8367.
1470 Simoes, L.C.; Simões, M. **2013**. Biofilms in drinking water: problems and solutions. *RSC advances*
1471 3(8), 2520-2533.
1472 Simpson, A.M.A.; Mitch, W.A. **2022**. Chlorine and ozone disinfection and disinfection byproducts
1473 in postharvest food processing facilities: A review. *Critical Reviews in Environmental Science*
1474 *and Technology* 52(11), 1825-1867.
1475 Sincak, P.; Ondo, J.; Kaposztasova, D.; Vircikova, M.; Vranayova, Z.; Sabol, J. **2014**. Artificial
1476 intelligence in public health prevention of legionellosis in drinking water systems. *International*
1477 *Journal of Environmental Research and Public Health* 11(8), 8597-8611.
1478 Singh, D.; Vardhan, M.; Sahu, R.; Chatterjee, D.; Chauhan, P.; Liu, S. **2023**. Machine-learning-
1479 and deep-learning-based streamflow prediction in a hilly catchment for future scenarios using
1480 CMIP6 GCM data. *Hydrology and Earth System Sciences* 27(5), 1047-1075.
1481 Singh, K.P.; Gupta, S. **2012**. Artificial intelligence based modeling for predicting the disinfection
1482 by-products in water. *Chemometrics and Intelligent Laboratory Systems* 114, 122-131.
1483 Siponen, S.; Jayaprakash, B.; Hokajärvi, A.-M.; Gomez-Alvarez, V.; Inkinen, J.; Ryzhikov, I.;
1484 Räsänen, P.; Ikonen, J.; Pursiainen, A.; Kauppinen, A. **2024**. Composition of active bacterial
1485 communities and presence of opportunistic pathogens in disinfected and non-disinfected
1486 drinking water distribution systems in Finland. *Water Research* 248, 120858.
1487 Sluban, B.; Lavrač, N. **2015**. Relating ensemble diversity and performance: A study in class noise
1488 detection. *Neurocomputing* 160, 120-131.
1489 Smith, A.; Sterba-Boatwright, B.; Mott, J. **2010**. Novel application of a statistical technique,
1490 Random Forests, in a bacterial source tracking study. *Water Research* 44(14), 4067-4076.
1491 Sokolova, M.; Lapalme, G. **2009**. A systematic analysis of performance measures for classification
1492 tasks. *Information Processing & Management* 45(4), 427-437.
1493 Speight, V.L.; Mounce, S.R.; Boxall, J.B. **2019**. Identification of the causes of drinking water
1494 discolouration from machine learning analysis of historical datasets. *Environmental Science:*
1495 *Water Research & Technology* 5(4), 747-755.
1496 Srivastav, A.L.; Patel, N.; Chaudhary, V.K. **2020**. Disinfection by-products in drinking water:
1497 Occurrence, toxicity and abatement. *Environmental Pollution* 267, 115474.
1498 Steel, R.G.D.; Torrie, J.H. **1960**. *Principles and procedures of statistics: with special reference to*
1499 *the biological sciences*. McGraw-Hill, New York.
1500 Stringer, C.; Wang, T.; Michaelos, M.; Pachitariu, M. **2021**. Cellpose: a generalist algorithm for
1501 cellular segmentation. *Nature Methods* 18(1), 100-106.
1502 Sudhakaran, S.; Amy, G.L. **2013**. QSAR models for oxidation of organic micropollutants in water
1503 based on ozone and hydroxyl radical rate constants and their chemical classification. *Water*
1504 *Research* 47(3), 1111-1122.

1505 Syafrudin, M.; Kristanti, R.A.; Yuniarto, A.; Hadibarata, T.; Rhee, J.; Al-Onazi, W.A.; Algarni, T.S.;
1506 Almarri, A.H.; Al-Mohaimed, A.M. **2021**. Pesticides in drinking water—a review.
1507 *International Journal of Environmental Research and Public Health* 18(2), 468.
1508 Taheran, M.; Naghdi, M.; Brar, S.K.; Verma, M.; Surampalli, R.Y. **2018**. Emerging contaminants:
1509 here today, there tomorrow! *Environmental nanotechnology, monitoring & management* 10,
1510 122-126.
1511 Tian, C.; Feng, C.; Chen, L.; Wang, Q. **2020**. Impact of water source mixture and population
1512 changes on the Al residue in megalopolitan drinking water. *Water Research* 186, 116335.
1513 U.S. EPA **2023** (accessed 2024). *Information about Public Water Systems*.
1514 <https://www.epa.gov/dwreginfo/information-about-public-water-systems>.
1515 Valbonesi, P.; Profita, M.; Vasumini, I.; Fabbri, E. **2021**. Contaminants of emerging concern in
1516 drinking water: Quality assessment by combining chemical and biological analysis. *Science of*
1517 *the Total Environment* 758, 143624.
1518 Walesch, S.; Birkelbach, J.; Jézéquel, G.; Haeckl, F.P.J.; Hegemann, J.D.; Hesterkamp, T.; Hirsch,
1519 A.K.H.; Hammann, P.; Müller, R. **2023**. Fighting antibiotic resistance—strategies and
1520 (pre)clinical developments to find new antibacterials. *EMBO Reports* 24(1), e56033.
1521 Wang, C.; Yang, H.; Liu, H.; Zhang, X.-X.; Ma, L. **2023**. Anthropogenic contributions to antibiotic
1522 resistance gene pollution in household drinking water revealed by machine-learning-based
1523 source-tracking. *Water Research* 246, 120682.
1524 Wen, X.; Chen, F.; Lin, Y.; Zhu, H.; Yuan, F.; Kuang, D.; Jia, Z.; Yuan, Z. **2020**. Microbial
1525 indicators and their use for monitoring drinking water quality—A review. *Sustainability* 12(6),
1526 2249.
1527 WHO **2011**. *Guidelines for Drinking-Water Quality*, 4th Edition. World Health Organization,
1528 Geneva.
1529 WHO **2021**. Antimicrobial resistance.
1530 Willard, J.; Jia, X.; Xu, S.; Steinbach, M.; Kumar, V. **2022**. Integrating scientific knowledge with
1531 machine learning for engineering and environmental systems. *ACM Computing Surveys* 55(4),
1532 66.
1533 Willmott, C.J.; Matsuura, K. **2005**. Advantages of the mean absolute error (MAE) over the root
1534 mean square error (RMSE) in assessing average model performance. *Climate Research* 30(1),
1535 79-82.
1536 Wingender, J.; Flemming, H.C. **2011**. Biofilms in drinking water and their role as reservoir for
1537 pathogens. *International Journal of Hygiene and Environmental Health* 214(6), 417-423.
1538 Wołos, A.; Koszelewski, D.; Roszak, R.; Szymkuć, S.; Moskal, M.; Ostaszewski, R.; Herrera, B.T.;
1539 Maier, J.M.; Brezicki, G.; Samuel, J.; Lummiss, J.A.M.; McQuade, D.T.; Rogers, L.;
1540 Grzybowski, B.A. **2022**. Computer-designed repurposing of chemical wastes into drugs.
1541 *Nature* 604(7907), 668-676.
1542 Wołowiec, M.; Komorowska-Kaufman, M.; Pruss, A.; Rzepa, G.; Bajda, T. **2019**. Removal of
1543 heavy metals and metalloids from water using drinking water treatment residuals as adsorbents:
1544 a review. *Minerals* 9(8), 487.
1545 Wolpert, D.H. **1992**. Stacked generalization. *Neural Networks* 5(2), 241-259.
1546 Wu, J.; Song, C.; Dubinsky, E.A.; Stewart, J.R. **2020**. Tracking major sources of water
1547 contamination using machine learning. *Frontiers in Microbiology* 11, 616692.
1548 Wu, Y.; Jiang, P.; Goh, S.G.; Yu, K.; Chen, Y.; He, Y.; Gin, K.Y.H. **2022**. Predicting relative risk of
1549 antimicrobial resistance using machine learning methods. *IFAC-PapersOnLine* 55(10), 1266-
1550 1271.

- 1551 Xiao, R.; Ou, T.; Ding, S.; Fang, C.; Xu, Z.; Chu, W. **2023**. Disinfection by-products as
1552 environmental contaminants of emerging concern: a review on their occurrence, fate and
1553 removal in the urban water cycle. *Critical Reviews in Environmental Science and Technology*
1554 53(1), 19-46.
- 1555 Xie, Y.; Sattari, K.; Zhang, C.; Lin, J. **2023**. Toward autonomous laboratories: Convergence of
1556 artificial intelligence and experimental automation. *Progress in Materials Science* 132, 101043.
- 1557 Xu, X.; Talbot, S.; Selvaraja, T. **2020**. ParasNet: Fast parasites detection with neural networks.
1558 *arXiv:2002.11327*, 10.48550/arXiv.42002.11327.
- 1559 Yaseen, Z.M. **2021**. An insight into machine learning models era in simulating soil, water bodies
1560 and adsorption heavy metals: Review, challenges and solutions. *Chemosphere* 277, 130126.
- 1561 Zainurin, S.N.; Wan Ismail, W.Z.; Mahamud, S.N.I.; Ismail, I.; Jamaludin, J.; Ariffin, K.N.Z.; Wan
1562 Ahmad Kamil, W.M. **2022**. Advancements in monitoring water quality based on various
1563 sensing methods: a systematic review. *International Journal of Environmental Research and*
1564 *Public Health* 19(21), 14080.
- 1565 Zanoni, M.G.; Majone, B.; Bellin, A. **2022**. A catchment-scale model of river water quality by
1566 Machine Learning. *Science of the Total Environment* 838, 156377.
- 1567 Zhang, C.; Brown, P.J.B.; Hu, Z. **2018**. Thermodynamic properties of an emerging chemical
1568 disinfectant, peracetic acid. *Science of the Total Environment* 621, 948-959.
- 1569 Zhang, C.; Brown, P.J.B.; Hu, Z. **2019a**. Higher functionality of bacterial plasmid DNA in water
1570 after peracetic acid disinfection compared with chlorination. *Science of the Total Environment*
1571 685, 419-427.
- 1572 Zhang, C.; Brown, P.J.B.; Miles, R.J.; White, T.A.; Grant, D.G.; Stalla, D.; Hu, Z. **2019b**. Inhibition
1573 of regrowth of planktonic and biofilm bacteria after peracetic acid disinfection. *Water Research*
1574 149, 640-649.
- 1575 Zhang, C.; Lu, J. **2021a**. Optimizing disinfectant residual dosage in engineered water systems to
1576 minimize the overall health risks of opportunistic pathogens and disinfection by-products.
1577 *Science of the Total Environment* 770, 145356.
- 1578 Zhang, C.; Lu, J. **2021b**. *Legionella*: a supplementary indicator of microbial water quality in
1579 municipal engineered water systems. *Frontiers in Environmental Science* 9, 684319.
- 1580 Zhang, C.; Qin, K.; Struewing, I.; Buse, H.; Santo Domingo, J.; Lytle, D.; Lu, J. **2021a**. The
1581 bacterial community diversity of bathroom hot tap water was significantly lower than that of
1582 cold tap and shower water. *Frontiers in Microbiology* 12, 625324.
- 1583 Zhang, C.; Struewing, I.; Mistry, J.H.; Wahman, D.G.; Pressman, J.; Lu, J. **2021b**. *Legionella* and
1584 other opportunistic pathogens in full-scale chloraminated municipal drinking water
1585 distribution systems. *Water Research* 205, 117571.
- 1586 Zhang, J.; Wang, Y.; Donarski, E.D.; Toma, T.T.; Miles, M.T.; Acton, S.T.; Gahlmann, A. **2022**.
1587 BCM3D 2.0: accurate segmentation of single bacterial cells in dense biofilms using
1588 computationally generated intermediate image representations. *npj Biofilms and Microbiomes*
1589 8(1), 99.
- 1590 Zhang, M.; Zhang, J.; Wang, Y.; Wang, J.; Achimovich, A.M.; Acton, S.T.; Gahlmann, A. **2020**.
1591 Non-invasive single-cell morphometry in living bacterial biofilms. *Nature Communications*
1592 11(1), 6151.
- 1593 Zhang, Y.; Gao, X.; Smith, K.; Inial, G.; Liu, S.; Conil, L.B.; Pan, B. **2019c**. Integrating water
1594 quality and operation into prediction of water production in drinking water treatment plants by
1595 genetic algorithm enhanced artificial neural network. *Water Research* 164, 114888.

1596 Zhong, S.; Zhang, K.; Bagheri, M.; Burken, J.G.; Gu, A.; Li, B.; Ma, X.; Marrone, B.L.; Ren, Z.J.;
1597 Schrier, J.; Shi, W.; Tan, H.; Wang, T.; Wang, X.; Wong, B.M.; Xiao, X.; Yu, X.; Zhu, J.-J.;
1598 Zhang, H. **2021**. Machine learning: new ideas and tools in environmental science and
1599 engineering. *Environmental Science & Technology* 55(19), 12741-12754.
1600 Zhou, Q.; Bian, Z.; Yang, D.; Fu, L. **2023a**. Stability of drinking water distribution systems and
1601 control of disinfection by-products. *Toxics* 11(7), 606.
1602 Zhou, Q.; Huang, J.; Guo, K.; Lou, Y.; Wang, H.; Zhou, R.; Tang, J.; Hou, P. **2023b**. Spatiotemporal
1603 distribution of opportunistic pathogens and microbial community in centralized rural drinking
1604 water: One year survey in China. *Environmental Research* 218, 115045.
1605 Zhou, X.; Tang, Z.; Xu, W.; Meng, F.; Chu, X.; Xin, K.; Fu, G. **2019**. Deep learning identifies
1606 accurate burst locations in water distribution networks. *Water Research* 166, 115058.
1607 Zhou, Z.; Zhong, D.; Zhang, Z.; Ma, W.; Chen, J.; Zhuang, M.; Li, F.; Zhang, J.; Zhu, Y.; Su, P.
1608 **2023c**. Biofilm on the pipeline wall is an important transmission route of resistome in drinking
1609 water distribution system. *Environmental Pollution* 335, 122311.
1610 Zhu, M.; Wang, J.; Yang, X.; Zhang, Y.; Zhang, L.; Ren, H.; Wu, B.; Ye, L. **2022**. A review of the
1611 application of machine learning in water quality evaluation. *Eco-Environment & Health* 1(2),
1612 107-116.
1613

Using Machine Learning to Ensure Drinking Water Quality

Water Source

- Source tracking
- Water source quality
- ...

Water Utility

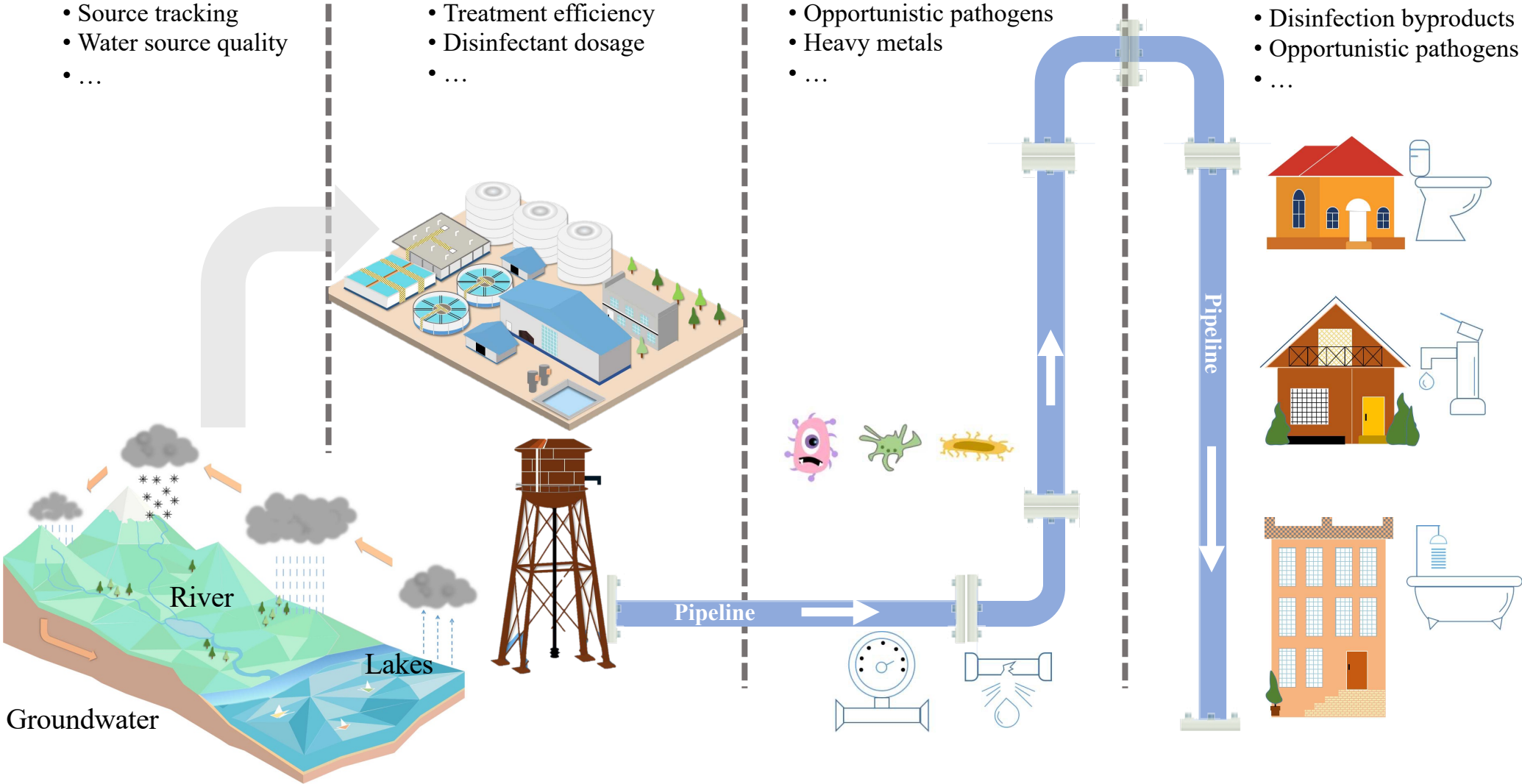
- Treatment efficiency
- Disinfectant dosage
- ...

Distribution System

- Opportunistic pathogens
- Heavy metals
- ...

Premise Plumbing

- Disinfection byproducts
- Opportunistic pathogens
- ...



Graphical abstract

Table 1 Machine learning to ensure safe drinking water supply from the physical perspective

| Topic | Task | Model | Inputs | Outputs | Metrics (Selected) | Performances (Selected) | Reference |
|--|------------------------------------|-------------------------------|--|--|--|---|------------------------|
| Predicting 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 prediction | 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 |
| Monitoring pipeline integrity | 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 Machine learning to ensure safe drinking water supply from the chemical perspective

| Topic | Task | Model | Inputs | Outputs | Metrics (Selected) | Performances (Selected) | Reference |
|--|---|--|---|--|--|--|----------------------|
| Optimizing drinking water disinfection | 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 | |
| Spectroscopic detection of DBPs | 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 |
| | | 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 |
| 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; 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;

Table 2 Machine learning to ensure safe drinking water supply from the chemical perspective

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 Machine learning to ensure safe drinking water supply from the chemical perspective (cont.)

| Topic | Task | Model | Inputs | Outputs | Metrics (Selected) | Performances (Selected) | Reference |
|--|---|-----------------------------|--|--|----------------------------------|--------------------------------|---------------------------------|
| Optimizing drinking water disinfection | 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 |
| Nitrification surveillance | Nitrification episodes classification | NB | 16S rRNA profiling | Nitrification episodes: stable or failure | AUC | 0.83 | Gomez-Alvarez and Revetta, 2020 |
| | Estimate NOx concentrations | SVR | NOx absorbances at various wavelengths | C_{NO_x} | RMSE and R^2 | RMSE < 0.04 | Hossain et al., 2021 |
| Heavy metal monitoring regulation | Pb ions concentration detection | SVR | S_{11} | Pb concentration | RMS | 0.71 | Oh et al., 2021 |
| | Spatial concentration mapping of heavy metal | 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 of Al residue | Kriging interpolation | Spatial and temporal data | Temporal-spatial distribution of residual Al | - | - | Tian et al., 2020 |
| | 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 and Mazloom, 2022 |
| | Heavy metal 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 |

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, SssO, sum of atom-type E-state; O^- , mean, mean intrinsic state values; <https://doi.org/10.26434/chemrxiv-2024-6541d-v2>, ORCID: <https://orcid.org/0009-0002-7289-9891> Content not peer-reviewed by ChemRxiv. License: CC BY 4.0

Table 2 Machine learning to ensure safe drinking water supply from the chemical perspective (*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; NB, naïve Bayes; AUC, area under the curve; NOx, nitrite and nitrate; 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.

Table 3 Machine learning to ensure safe drinking water supply from the microbiological perspective

| Topic | Task | Model | Inputs | Outputs | Metrics (Selected) | Performances (Selected) | Reference |
|--|--|------------------------------|---|---|-----------------------------------|------------------------------------|---|
| Surveilling and mitigating 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 higher-risk, serogroup and contamination levels prediction of <i>Legionella</i> spread | SaTScan, XGB, LR, and SVM | Survey, spatial and meteorological info., and risk level to <i>Legionella</i> ; | Higher-risk level clusters; serogroup of a sample and the contamination level | Accuracy and F1-score | XGBoost > SVM > LR | Brunello et al., 2022 |
| Analyzing drinking water microbial communities | 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 |
| | RF | ARA | Source classification | ARCC | 82.3% | Smith et al., 2010 | |
| | Microbial contamination prediction | 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 | - | Pinto et al., 2014 |
| | | UniFrac | - | - | Unweighted/weighted UniFrac score | - | Lozupone et al., 2011; Bruno et al., 2018; Ling et al., 2018; Li et al., 2017 |

NARA, Neural network designed on approximate reasoning architecture; Q , flow rate; 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 3 Machine learning to ensure safe drinking water supply from the microbiological perspective (cont.)

| Topic | Task | Model | Inputs | Outputs | Metrics (Selected) | Performances (Selected) | Reference |
|---|--|--------------------------|---|---|---|---|-----------------------------|
| Parasite detection | 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 |
| | | | | Multiple classification or binary classification | Accuracy, precision, recall, and F1-score | Accuracy > 99.6% | Luo et al., 2021 |
| | <i>Cryptosporidium</i> and <i>Giardia</i> contamination intensity prediction | LDFA | Microbiological, physicochemical, and meteorological parameters | (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 |
| Biofilm development assessment | Biofilm development analysis | RT and RF | System physical and hydraulic characteristics, sampling and incubation, and physico-chemical of water | HPC | <i>R</i> | 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 |
| Risk analysis and source tracking of antimicrobial resistance | 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 pollution 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 | <i>r</i> | <i>r</i> = 0.98 | Wang et al., 2023 |

LDFA, linear discriminant function analysis; RT, regression trees; HPC, heterotrophic plate count; OSA, over-segmentation abundances; IoU, intersection-over-union; 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; ARG, Antimicrobial resistance genes

Table 4 Machine learning to ensure safe drinking water supply from the temporal perspective

| Topic | Task | Model | Inputs | Outputs | Metrics (Selected) | Performances (Selected) | Reference |
|--|---|---------------------------------------|---|----------------------------------|------------------------------------|-------------------------|------------------------------|
| Anomalies and contamination events detection | Anomalies 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.

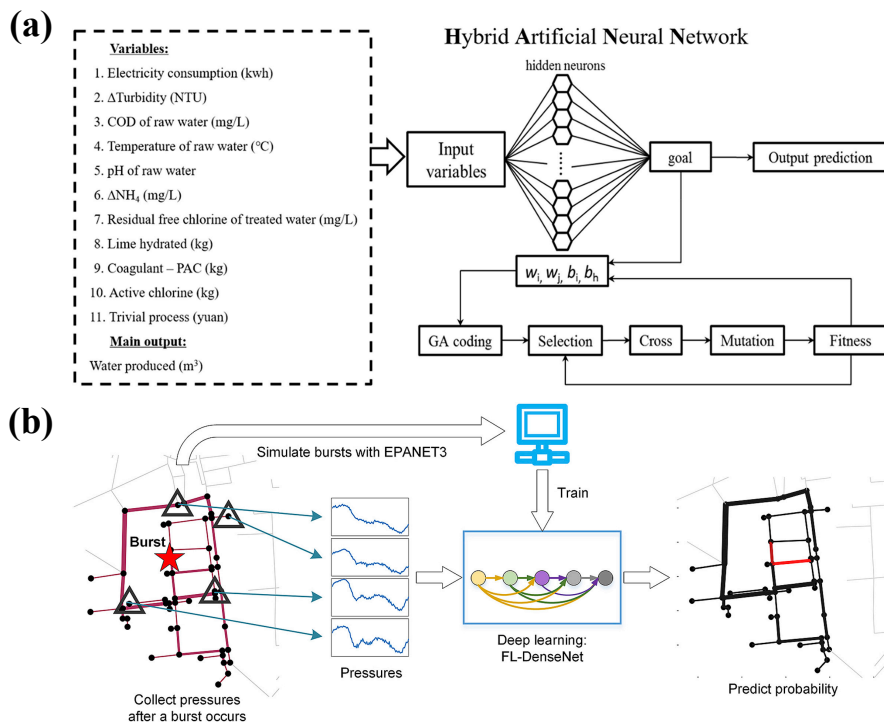


Figure 1. (a) Input and output variables used for modeling and the proposed hybrid artificial neural network framework for prediction of drinking water production. Reproduced with permission from Zhang et al., 2019. Copyright 2019 Elsevier. (b) Schematic of fully-linear DenseNet (BLIFF) model for accurate identification of burst locations in EWS networks. Reproduced with permission from Zhou et al., 2019 (CC BY 4.0).

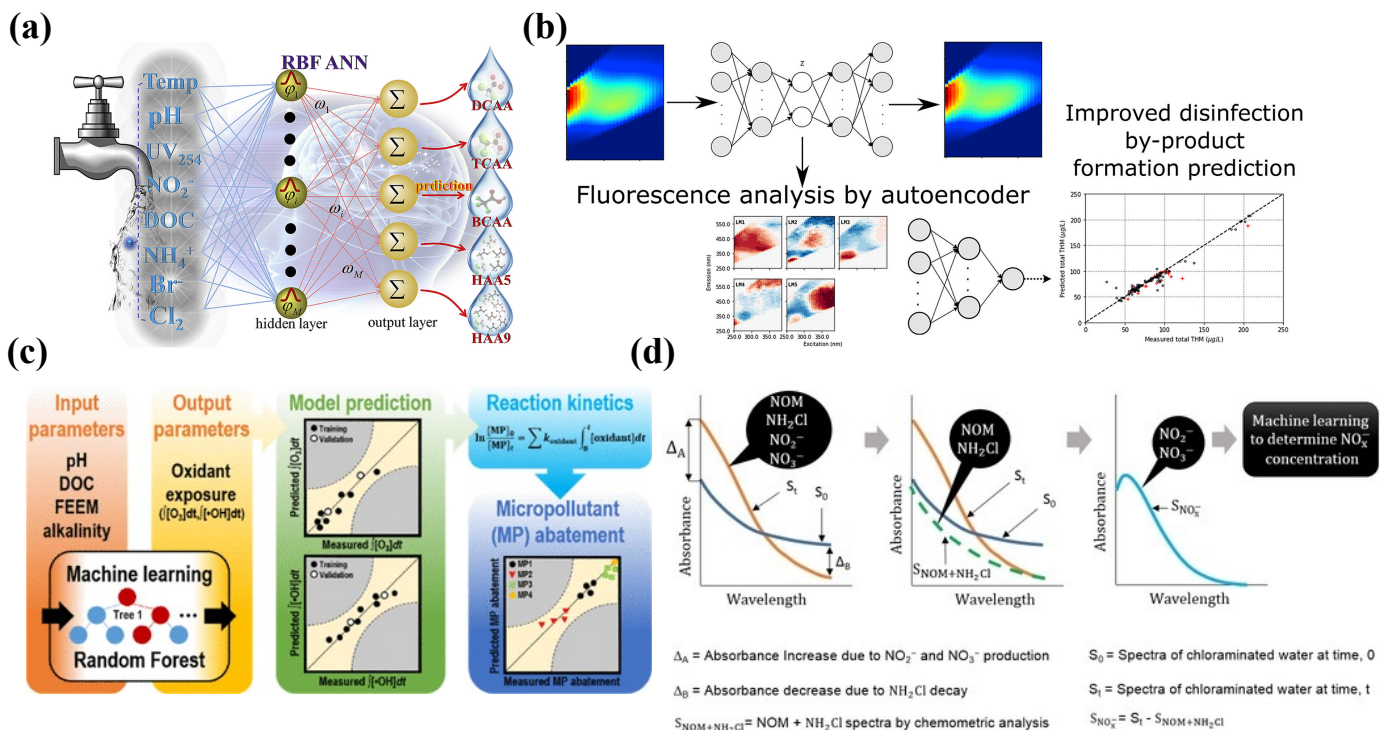


Figure 2. (a) Schematic of radial basis function (RBF) artificial neural network (ANN) model for prediction of disinfection by-products (DBPs). Reproduced with permission from Lin et al., 2020, Copyright 2020 Elsevier. (b) Schematic of autoencoder model for prediction of DBPs. Reproduced with permission from Peleato et al., 2018, Copyright 2018 Elsevier. (c) Schematic of random forest (RF) model for prediction of micropollutant abatement. Reproduced with permission from Cha et al., 2021, Copyright 2021 American Chemical Society. (d) Prediction of nitrate and nitrite concentrations over support vector regression (SVR) model. Reproduced with permission from Hossain et al., 2021 (CC BY 4.0).

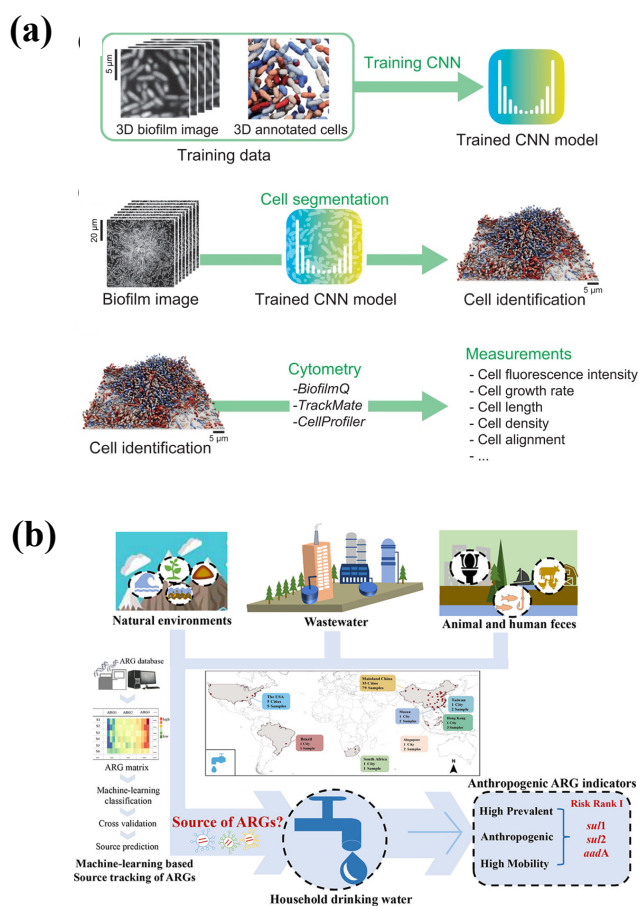


Figure 3. (a) Deep-learning-based workflow for single-cell measurements in three-dimensional biofilms. Reproduced with permission from Jelli et al., 2023, Copyright 2023 Elsevier. (b) SourceTracker was performed to investigate the pollution sources of antimicrobial resistance genes (ARGs) in household drinking water. Reproduced with permission from Wang et al., 2023, Copyright 2023 Elsevier.

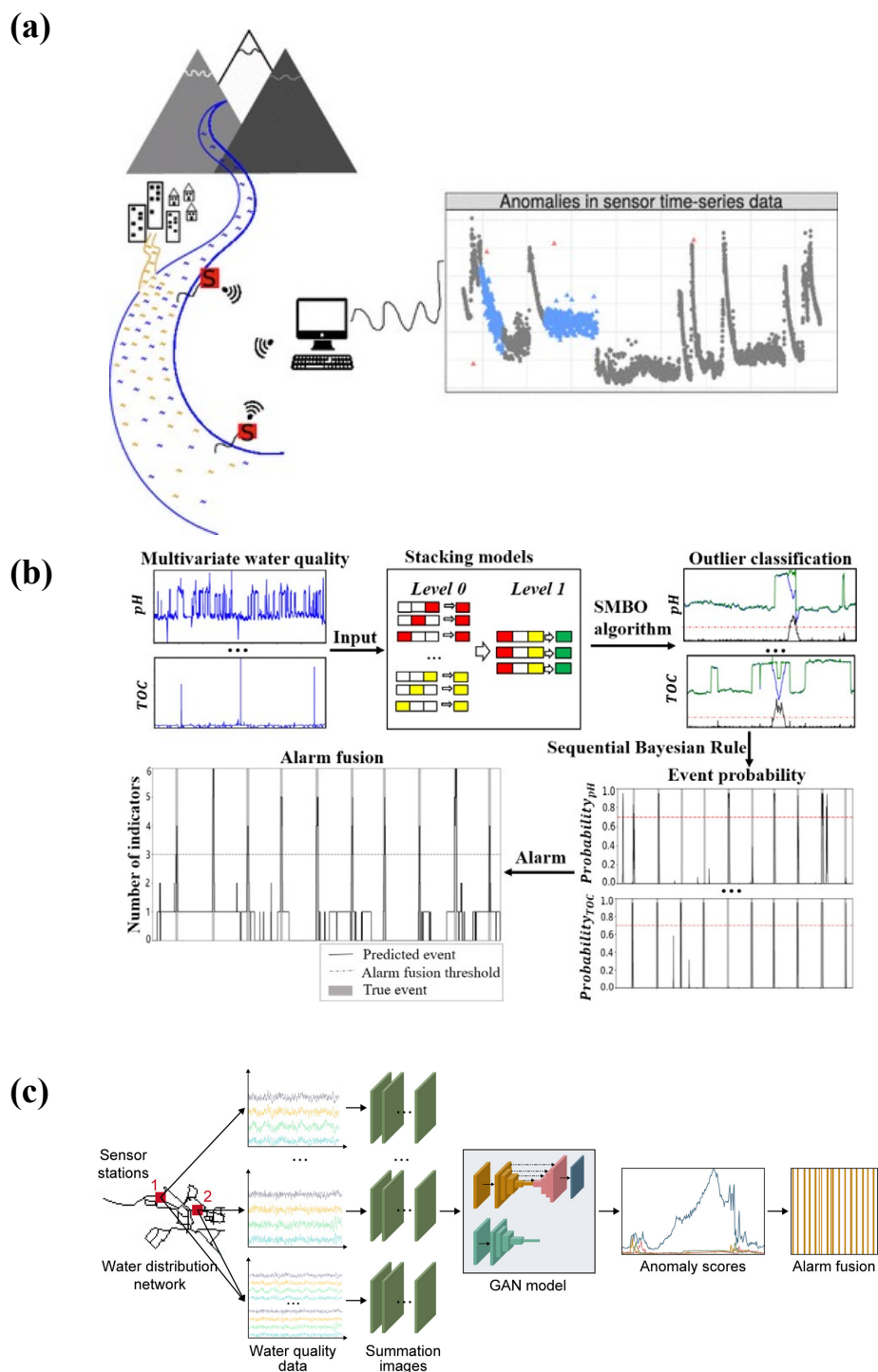


Figure 4. (a) Detection of technical anomalies in water quality using artificial neural network (ANN) model. Reproduced with permission from Rodriguez-Perez et al., 2020, Copyright 2020 American Chemical Society. (b) A stacking ensemble model for contamination event detection using multiple water quality parameters. Reproduced with permission from Li et al., 2022, Copyright 2022 Elsevier. (c) Detection of contamination events using generative adversarial network (GAN) model. Reproduced with permission from Li et al., 2023, Copyright 2023 Elsevier.

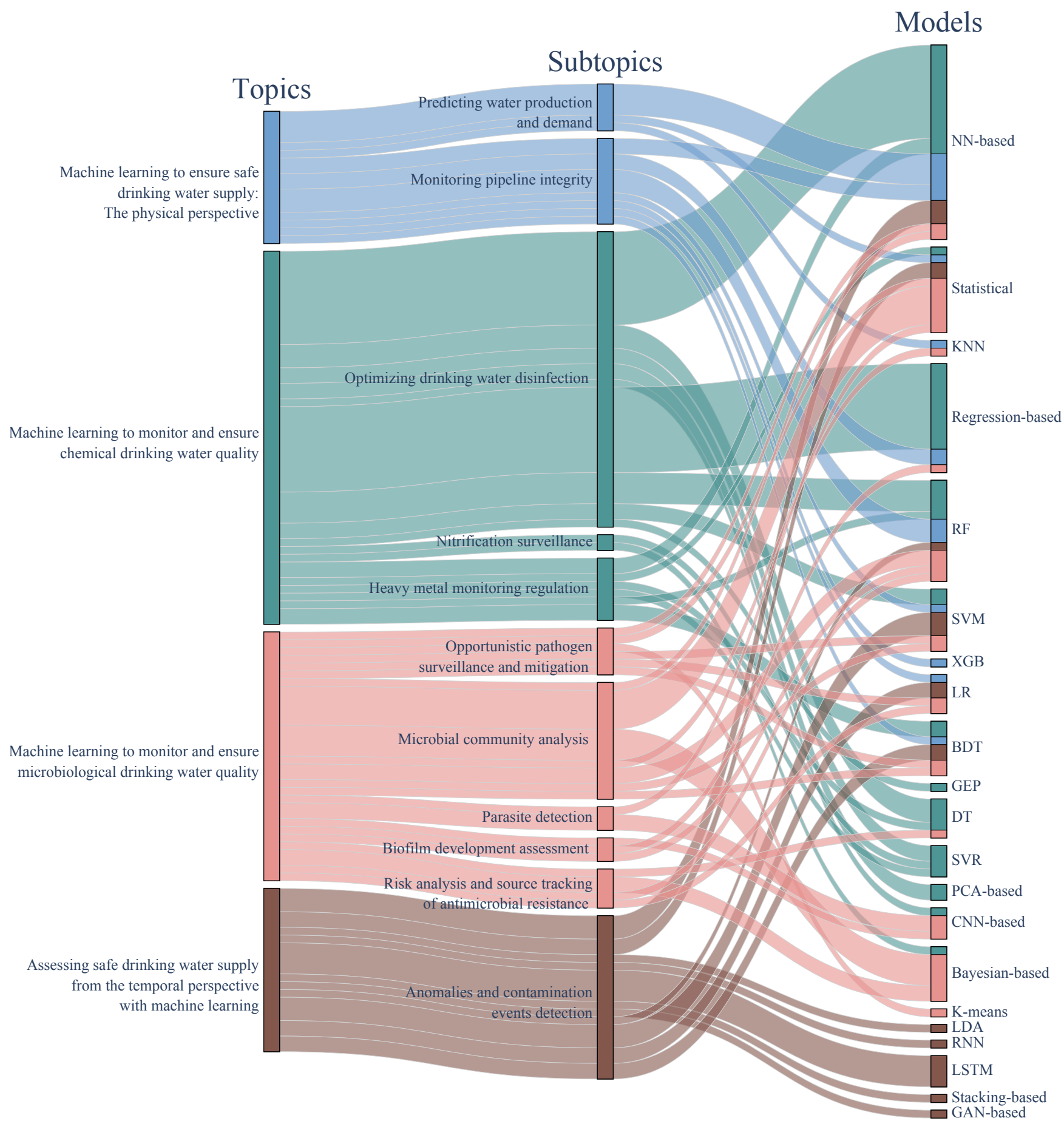


Figure 5. Machine learning model distribution across drinking water supply research topics

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 neutral network; LDA, linear discriminant analysis; RNN, recurrent neural network; LSTM, long short-term memory; GAN, generative adversarial network.