Machine learning-based surrogate modelling for Urban Water Networks: Review and future research directions

Alexander Garzón^{1,1}, Zoran Kapelan^{2,2}, jeroen langeveld^{3,3}, and Riccardo Taormina^{3,3}

¹Delft University of Technology ²TU Delft, Faculty of Civil Engineering and Geosciences ³TU Delft

November 30, 2022

Abstract

Surrogate models replace computationally expensive simulations of physically-based models to obtain accurate results at a fraction of the time. These surrogate models, also known as metamodels, have been employed for analysis, control, and optimisation of water distribution and urban drainage systems. With the advent of machine learning (ML), water engineers have increasingly resorted to these data-driven techniques to develop metamodels of urban water networks. In this manuscript, we review 31 recent papers on ML-based metamodeling of urban water networks to outline the state-of-the-art of the field, identify outstanding gaps, and propose future research directions. For each paper, we critically examined the purpose of the metamodel, the metamodel characteristics, and the applied case study. The review shows that current metamodels suffer several drawbacks, including i) the curse of dimensionality, hindering implementation for large case studies; ii) black-box deterministic nature, limiting explainability and applicability; and iii) rigid architecture, preventing generalization across multiple case studies. We argue that researchers should tackle these issues by resorting to recent advancements in ML concerning inductive biases, robustness, and transferability. The recently developed Graph Neural Network architecture, which extends deep learning methods to graph data structures, is a preferred candidate for advancing surrogate modelling in urban water networks. Furthermore, we foresee increasing efforts for complex applications where metamodels may play a fundamental role, such as uncertainty analysis and multi-objective optimisation. Lastly, the development and comparison of ML-based metamodel can benefit from the availability of new benchmark datasets for urban drainage systems and realistic complex networks.

1 2	manuscript submitted to Water Resources Research Machine learning-based surrogate modelling for Urban Water Networks: Review and future research directions
3	A. Garzón ¹ , Z. Kapelan ¹ , J. Langeveld ^{1,2} and R. Taormina ¹
4 5 6	¹ Faculty of Civil Engineering and Geosciences, Department of Water Management, Delft University of Technology, Stevinweg 1, 2628 CN Delft, The Netherlands ² Partners4UrbanWater, Graafseweg 274, 6532 ZV Nijmegen, The Netherlands
7 8	Corresponding author: Alexander Garzón (J.A.GarzonDiaz@tudelft.nl)
9	Key Points:
10 11	• Machine Learning surrogate models have been widely employed for a variety of applications concerning urban water networks.
12 13	• New research should focus on machine learning metamodels that account for inductive biases, robustness, and transferability.
14 15	• Further research should focus on complex problems involving uncertainty and multi-objective optimisation, as well as improved benchmarking.

16 Abstract

Surrogate models replace computationally expensive simulations of physically-based models to obtain 17 accurate results at a fraction of the time. These surrogate models, also known as metamodels, have been 18 employed for analysis, control, and optimisation of water distribution and urban drainage systems. With 19 the advent of machine learning (ML), water engineers have increasingly resorted to these data-driven 20 techniques to develop metamodels of urban water networks. In this manuscript, we review 31 recent papers 21 on ML-based metamodeling of urban water networks to outline the state-of-the-art of the field, identify 22 outstanding gaps, and propose future research directions. For each paper, we critically examined the 23 purpose of the metamodel, the metamodel characteristics, and the applied case study. The review shows 24 that current metamodels suffer several drawbacks, including i) the curse of dimensionality, hindering 25 implementation for large case studies; ii) black-box deterministic nature, limiting explainability and 26 applicability; and iii) rigid architecture, preventing generalization across multiple case studies. We argue 27 that researchers should tackle these issues by resorting to recent advancements in ML concerning inductive 28 biases, robustness, and transferability. The recently developed Graph Neural Network architecture, which 29 extends deep learning methods to graph data structures, is a preferred candidate for advancing surrogate 30 modelling in urban water networks. Furthermore, we foresee increasing efforts for complex applications 31 where metamodels may play a fundamental role, such as uncertainty analysis and multi-objective 32 optimisation. Lastly, the development and comparison of ML-based metamodel can benefit from the 33 availability of new benchmark datasets for urban drainage systems and realistic complex networks. 34

35 Plain Language Summary

Analysis and improvement of urban water networks requires hydrodynamic models. Since these models 36 are computationally expensive, researchers and engineers often resort to fast alternatives known as 37 surrogate models. With the rise of artificial intelligence, machine learning methods have been increasingly 38 used for surrogate modelling of urban water networks. In this study, we thoroughly reviewed recent papers 39 on the field to outline the current state-of-the-art and propose future research directions. While many 40 successful applications already exist, we found that these models have three main limiting factors: i) they 41 need large amounts of data, ii) they are not explainable, and iii) they are too specific to each case. We argue 42 that researchers can overcome these limitations by considering recent advancements in artificial intelligence 43 and implement modeling techniques that better leverage the structure of the underlying data. Other 44 promising direction include developing comprehensive benchmark databases and leveraging surrogate 45 models for more complex applications. 46

47 **1 Introduction**

61

Urban water networks (UWNs) comprise drinking water distribution and urban drainage systems 48 (WDS and UDS). The former are responsible for supplying drinking water to cities and the latter for 49 evacuating wastewater and stormwater runoff. These infrastructures are a fundamental part of the city and 50 are directly linked to its development (Brown et al., 2009). Each of these systems faces challenges to 51 improve and maintain quality service in a dynamic urban environment under a widening range of climatic 52 conditions; especially, in a climate-changing situation. Designing, optimising, and intervening in these 53 systems requires approximating their hydraulic behaviour. Several models have been developed in the past 54 years for simulating UWNs. Traditional modelling approaches are either based on accurate description of 55 the physical processes or rely on simplified conceptual approaches; nonetheless, the former usually entail 56 computationally expensive calculations while the latter lack fidelity. Applications such as optimisation, 57 real-time modelling, and uncertainty analysis need an efficient model for evaluating the performance of a 58 system multiple times or as fast as possible. Consequently, they require short execution times while 59 maintaining a sufficient level of detail. 60

1.1 Surrogate modelling

Water modellers have resorted to surrogate models (SMs) to replace computationally costly models. Following the classification given by Razavi et al. (2012b), SMs, also known as metamodels or reducedorder models, can be categorized as Lower-fidelity Physically-based surrogates (LFPB) or response surface (RS) surrogates. On one hand, LFPB metamodels modify the original model to reduce its computational

effort. These models simplify the original model by lowering the resolution (e.g., larger time-steps) of the 66 output or replacing computationally costly components with faster alternatives or complements (e.g., 67 kriging, linear regression, neural networks (Fernandez et al., 2017)). On the other hand, RS surrogates avoid 68 using the original model and replace it altogether with a faster-to-run alternative. In this branch of SMs, the 69 original model is perceived as an input-output function and the metamodel is used to mimic the output 70 surface as best as possible. Some of the algorithms for approximating response surfaces are polynomial 71 interpolation, kriging, and more recently, machine learning (ML) algorithms. The following paragraphs 72 summarize the advantages and disadvantages of LFPB and RS metamodels according to Razavi et al. 73 74 (2012b).

75 Lower-fidelity Physically-based surrogates (LFPB), also known as multifidelity based surrogates or "coarse" models, include techniques such as network simplification (Dempsey et al., 1997; 76 Paluszczyszyn et al., 2013; Ulanicki et al., 1996), and skeletonization (Shamir et al., 2008). Compared 77 against RS metamodels, LFPB surrogates are expected to better emulate the unexplored regions of the 78 explanatory variable (input) space (i.e., regions far from the previously evaluated points with the high-79 fidelity model) and, as such, perform more reliably in extrapolation. As for their drawbacks, LFPB models 80 81 rely on the assumption that high-fidelity and low-fidelity models share the basic features and are correlated in some way. If this assumption is not satisfied, the surrogate modelling framework would not work, or 82 provide minimal gains. Moreover, mapping the outputs from low resolution to the original resolution is not 83 a trivial task, and may add complexity or uncertainty to the estimations. 84

Response surface (RS) surrogates, also known as statistical and black-box models, include 85 techniques such as polynomials (Schultz et al., 2004), kriging (Baú & Mayer, 2006), and neural networks 86 (Behzadian et al., 2009). Some of their advantages include the possibility of maintaining the fidelity of the 87 original model, being model-independent (i.e., not requiring access to the components, such as code or 88 equations of the original model), and easier implementation with respect to LFPB surrogates. Nonetheless, 89 they can be hard to train for high-dimensional problems, which may require extreme computational costs 90 to create large enough databases to train the metamodels. Moreover, RS metamodels require scrupulous 91 validation to minimize the chance of over-fitting and maximize their ability to extrapolate. 92

93 **1.2 Machine learning methods**

ML methods are part of artificial intelligence (AI) which is a broad term for tools that mimic 94 cognitive human capabilities. The use of AI has rapidly increased in recent years. The number of peer-95 reviewed publications across all fields between 2000 and 2019 has grown around 12 times (D. Zhang et al., 96 97 2021) and with them, multiple algorithms, architectures, and tools have been created. Fields in which ML methods have shown outstanding results include computer vision, speech recognition, and language 98 processing. Most of these applications use supervised learning, which identifies a branch of ML that is 99 similar to RS metamodelling. Supervised ML employs a set of input-output examples, also known as the 100 labelled training dataset, to calibrate a model by minimizing the error between the model predictions and 101 the values assumed as ground truth. This set of algorithms usually increase their performance at a given 102 task as the amount of labelled examples grows larger. Due to their successes, supervised ML methods, and 103 in particular deep learning (DL) and artificial neural networks (ANNs), are widely employed for surrogate 104 modelling across many fields of science and engineering (Liu et al., 2021; Peng et al., 2020; Wu et al., 105 106 2020). Although scientific studies on ML applications for water resources date back to over two decades ago (Maier & Dandy, 2000), Hadjimichael et al. (2016) noted that this trend is not necessarily witnessed in 107 the urban water sector. 108

109

1.3 Previous studies - Surrogate Modelling in Urban Water Networks

Previous studies have reviewed the application of metamodels in water resources. Razavi et al. (2012b) outline taxonomies, practical details, and advances of these SMs in water resources along with recommendations for future research. Among the multiple insights of this work, they highlight the nontrivial effort to choose the right metamodel approach to the problem at hand and advocate for further research on these methods, especially in their assessment and validation. Furthermore, in the same year, Razavi et al. (2012a) numerically assessed metamodeling strategies in computationally intensive optimization, showing that metamodeling is not always a reliable approach, especially for complex

response surfaces. The authors also warned about the inappropriateness of neural network models when having a limited computational budget. Later, Broad et al., (2015) presented a formalized qualitative process to determine the most suitable scope for a metamodel based on the evaluation of a fitness function to maximize fidelity. Hadjimichael et al. (2016) reviewed the application of AI methods to UWS management and their integration with decision support systems. While valuable, these published reviews give low emphasis to SMs for UWNs, and do not account for the recent growth in machine learning-based surrogate models (MLSMs) driven by the rapid advancements in AI.

This study aims to fill this gap by assessing the current state of MLSMs for UWNs in order to propose future directions based on identified outstanding issues and recent developments in ML. To achieve this purpose, we applied the review methodology described in Section 2 to review 31 published applications of metamodels for water networks. The results of the review are reported and discussed in Section 3, while major current gaps are detailed in Section 4. We propose future research directions in Section 5 and provide conclusions in Section 6.

130 2 Materials and Methods

We conducted a semi-systematic (Snyder, 2019) review of MLSM applications for UWNs to synthesize the state-of-the-art of the field. The review integrates the multiple applications of metamodels across water network applications, and explores them in a transversal manner. First, we searched journal papers in which MLSMs were applied to UWNs. Second, we determined a set of criteria to assess the relevant characteristics when applying these metamodels to UWNs' problems.

136 **2.1 Search methodology**

We reviewed journal papers published in the last two decades (2001-2021) that use MLSMs for 137 WDSs and UDSs. We established two main search criteria: surrogate modelling and water networks. Since 138 both topics have a multiplicity of names, each of them was represented by a set of keywords. For surrogate 139 modelling, the search terms were: "Surrogate model*", "Metamodel*", "Response surface", "model 140 emulation", and "hybrid model". In the case of water networks, the search terms referred to both water 141 distribution and drainage systems along with popular software for their analysis, "Water distribution", 142 "Water supply", "Drinking water", "Urban drainage", "Wastewater", "Sewer", "Sewerage", "EPANET", 143 "WaterCAD", "SWMM", "SOBEK", and "Urban water". 144

For the search, we employed the SCOPUS database. By intersecting the search terms, we identified an initial set of 64 papers that were further filtered to only include ML applications, yielding a total of 31 papers to review. Next, we searched through the citations of the selected set of papers and other relevant papers in the field (i.e., Maier et al., 2014; Maier & Dandy, 2000; Razavi et al., 2012b) for further references. However, the original set already contained the cited papers. Therefore, the results are equivalent to the keyword search. This validates the thoroughness of the original search and makes the methodology more replicable by avoiding arbitrarily selected papers.

This list of papers may not be totally inclusive since some studies do not use the formal terminology of surrogate modelling, as indicated by Razavi et al. (2012b). Nevertheless, the purpose of this paper is to depict the recent state-of-the-art, identify gaps in knowledge and propose future research directions. We believe that the selected set of papers is sufficient to achieve this goal.

156 **2.2 Analytical methodology**

In addition to the search criteria, it was necessary to establish an analytical framework that allowed to classify, compare, and evaluate the application of the metamodels across the collected literature. To

achieve this, we identified the most relevant aspects of each paper in three broad categories: i) purpose, ii)case study, and iii) metamodel.

Purpose includes general information about the application of the metamodel. It includes the type of network (distribution or drainage) and the application category (e.g., optimisation, real-time) as major grouping categories. In addition, it includes the specific application (e.g., optimisation of operation, realtime for flood prediction) as a more detailed description for each paper.

Case study contains information on the physical water network used for the testing and validation of a developed metamodel. This includes the name or location of the case study, whether it is a real case or a benchmark, and its size, indicated by the number of pipes or by the area. The size attribute is also reported as a categorical value ranging from small (S) to large (L), as shown in Table 1.

Size	Number of pipes in the simulation model	Area [km ²]
Small (S)	<100	<5
Medium (M)	101-250	5 - 10
Intermediate (I)	251-500	10 - 20
Large (L)	>500	>20

169 **Table 1.** Categories of network size based on number of pipes or area

Metamodel reports details on the computational algorithm (e.g., ANNs, Support Vector Machines) 170 used to replace the original simulator along with further details on its architecture (i.e., deviations from a 171 hidden layer ANN). The type and number of input and output variables are also reported to infer the 172 dimensionality of the SM and the complexity of the RS to approximate. As for the performance, we report 173 the computational speed-up provided by the metamodel and the fidelity to the original simulation, usually 174 approximated with an accuracy metric. These criteria have been identified as the most relevant ones by 175 previous related studies (Broad et al., 2015; Razavi et al., 2012b). Nevertheless, it is possible to consider 176 177 other factors, such as development time, robustness and explainability. While assessing these criteria may enrich the analysis, they are not employed in most of the surveyed papers, and they are thus not included 178 in this review. 179

180 **3 Review – Current status of Machine Learning Surrogate Models in Urban Water networks**

The analysis of the surveyed papers show an increase in research activity between 2015 and 2020 181 with approximately two-thirds of the manuscripts published during this period. In terms of application, 182 most of these papers are related to optimisation. For the case study, there is a noticeable difference between 183 WDSs and UDSs since the latter networks lack the use of benchmark cases. Regarding the metamodel, the 184 most popular algorithm is the fully connected ANN; because of this, we report the details of the used 185 metamodel as deviations from a standard, one hidden layer, fully connected ANN, also referred to as simple 186 Multi-layer perceptron (MLP). Table 2 summarizes the extracted information of the reviewed papers 187 arranged in the previously mentioned categories: purpose, case study, and metamodel. 188

Table 2. *List of reviewed papers and metamodeling approaches.*

Purpose			Case study				Metamodel			Metamodel Performance		
Water network	Application category	Reference	Application	Location	Size: Pipes in model / [area km^2]	Classification by size	Туре	Deviations from simple MLP	Inputs (Number)	Outputs (Number)	Computational saving	Accuracy
		(Sayers et al., 2019)	Design	TLN, GOY, MOD. BIN	8, 30, 317, 454	S, S, I, I	Benchmark	2 hidden layers	Diameters *	Rating of the network (1)	Not reported	Not reported
	Optimisation	(Dini & Tabesh, 2019)	Renovation planning	TLN and Ahar, Azerbaijan	8 and 192	S, M	Benchmark and Real case		Diameters *	Nodal pressure* and chlorine concentration *	Not reported	Not reported
		(Dini & Tabesh, 2017)	Model calibration	TLN and Ahar, Azerbaijan	8 and 192	S, M	Benchmark and Real case		Observed residual chlorine *	Wall Decay coefficient (1)	58x faster (98.3%)	Average error (3.85%)
		(Andrade et al., 2016)	Design	HAN and Maricopa, Arizona	34 and 1090	S, L	Benchmark and Real case	Comparison of ANNs varying number of inputs and outputs	Diameters and Chlorine dosing rates	Chlorine concentration. (HAN): 3; (Maricopa): 9	Not reported	NSE (~90%)
		(Bi & Dandy, 2014)	Design	(I) NYT, (II) modified NYT and (III) Jilin	21, 21, and 34	S, S, S	(I) Benchmark, (II) modified benchmark, and (III) synthetic network		Diameters and Chlorine dosing rates (I & II: 22; III: 35)	Pressures at some nodes (I & II: 4; III: 5) and residual chlorine at one node (I & II: 1; III: 7)	(I & II) 91%; (III) 93%, 88%, and 77%	MSE (Not reported, 0.001 as one stopping criteria)
		(Broad et al., 2010)	Operation	Wallan, Australia	2097;(Sk: 1376)	L (L)	Real case		Trigger levels (45) and Chlorine rates (5)	Pressure Head at critical node (1), Chlorine residual (1), energy value (1), or Total chlorine dosed (1)	99%	NSE (~0.6 for the full model, ~0.98 for skeletonized model)
		(Behzadian et al., 2009)	Sensor placement	Anytown; Mahalat, Iran	41, and 1814;(Sk: 217)	S, L (M)	Benchmark and Real case		Available sensors	Sampling design accuracy (1)	8x and 25x faster (87% and 96%)	Pareto similarity: 93%
Water		(Salomons et al., 2007)	Operation	Haifa-A, Israel	126	м	Modified real case		Pumping status (13), Valve settings (1), DMA demands (6), Storage levels (9)	Power consumption (5), pressures (4), future storage levels (9)	25x faster (96%)	RMSE (0.481%) ~5 cm averaged over all tanks
distribution systems		(Martínez et al., 2007)	Operation	Valencia, Spain	772	L	Modified real case		Pumping status (6), Valve settings (10), DMA demands (6), Storage levels (2)	Power consumption (6), flow rates (3), pressures (4), future storage levels (2)	94x faster (99%)	RMSE (1.30%)
		(Broad et al., 2005a)	Design	NYT	21	S	Benchmark		Diameters and Chlorine dosing rate (22)	Four pressure nodes (1) or Chlorine concentration (1)	700x faster (99.85%)	RMSE (0.05 - 0.250)
	Real-time	(Pasha & Lansey, 2014)	Warm solutions for pump scheduling	Modified Anytown	37	S	Modified Benchmark	SVM	Pump combination, demand multiplier, initial tank levels	Energy and final tank levels	84.25%	NSE (0.99)
		(Rao & Alvarruiz, 2007; Rao & Salomons, 2007)	Real-time pump scheduling	Modified AnyTown	41	s	Modified Benchmark		Number of operating pumps (1), aggregated demand (1), and tank levels (3)	Power consumption (1), pressures (3), new tank levels (3)	10-fold (90%)	RMSE (1.65%)
	Uncertainty analysis	(Yoon et al., 2020)	Seismic risk assessment	A-city, South Korea	85	S	Anonymous real case	15 layers - Deep neural network	Components' state (218)	Network performance (1)	99%	<5%
		(Beh et al., 2017)	Planning under deep uncertainty	Adelaide, Australia	NA	L	Real case	Combination of 4 MLPs	Supply augmentation options (9) and Uncertain variables: Population and climate change scenarios (2)	(I) PV of cost (II) PV of Greenhouse gases (III) Reliability (IV) Vulnerability	>99%	Relative error (+-5%) NSE (~0.94, 0.95, 0.78, and 0.84)
	System state estimation	(Lima et al., 2018)	Nodal pressure estimation at near real-time	Campos do Conde II and Cambuí, Brazil	153 and 167	М, М	Real case		Pressure in sensors Steady State: (3) - Extended (24h): 96. Cambuí: (4)	Pressure in nodes Steady State: (118) - Extended (24h): 2832. Cambuí: Steady (154 and 4)	Not reported	Relative error (<1%) and (<4%)
		(Meirelles et al., 2017)	Calibration with estimated pressures	Campos do Conde II, Brazil and C- Town	153 and 429	M, I	Real case and Benchmark		Pressure in sensors Steady State: (3) - Extended (24h): 96. C-Town: 5 MLPs, one per DMA.	Pressure in nodes Steady State: (118) - Extended (24h): 2832	Not reported	Average error (0.15 m) Max. Error (13.83 m)

Purpose			Case study			Metamodel			Metamodel performance			
Water network	Application category	Reference	Application	Location	Size # Pipes in model /[area km ²]	Classification by size	Туре	Deviations from simple MLP	Inputs (Number)	Outputs (Number)	Computational saving	Accuracy
Urban drainage systems	Optimisation	(Seyedashraf et al., 2021)	Design	Bogotá, Colombia; Windsor, Canada	511 and 122	L, M	Stormwater - Real cases	Generalized regression - 2 hidden layers	SUDS characteristics: area, type, and location (20)	Boundary condition: Inflow (1)	95%	Mean error (<0.015) CC (0.99)
		(W. Zhang et al., 2019)	Design	Urban catchment in China	182	М	Stormwater* - Real case	Ensemble of 100 MLPs	Tank length and width (2)	Flood depth (1) or peak flow (1)	80 - 90 %	NSE (Between 0.66 and 0.92 depending on the rainfall scenario)
		(Raei et al., 2019)	Design	Tehran, Iran	[20 km²]	I	Stormwater* - Real case	2 hidden layers	Area sizes of the LID, Imperviousness and rainfall (3), TSS/BOD build-up (+1), TSS/BOD wash-off (+1)	The volume of runoff (1) or BOD (1) or TSS (1)	Not reported	NSE (0.99)
		(Latifi et al., 2019)	Design	Tehran, Iran	[20 km ²]	I	Stormwater* - Real case		Rainfall value, 6 build-up coefficients, 6 wash off coefficients, 6 imperviousness coefficients, and 32 values for area and type of LIDs (51)	Runoff volume, BOD, TSS (3)	Not reported	Not mentioned
		(Huang et al., 2015)	Design	Zhong-He district, Taiwan	[20.29 km ²]	L	Stormwater* - Real case		Catchment precipitation, Full pipe percentage of water flow in 3 points, the quantity and capacity of rain barrels in four regions (12)	Water level/flooding at t + 1 (1)	Not reported	MAE (<15%) CC (>0.94 ~0.97)
	Real-time	(Kim & Han, 2020)	Flood prediction	Seoul, Korea	[3.19 km ² *]	М	Stormwater* - Real case	8 hidden layers	Total rainfall, Max. Rainfall in 1 - 3 hours, rainfall intensity, statistics (SD, Skewness, kurtosis), inter-event time (9)	Total accumulative overflow (1)	~99%	Mean relative errors between 2% - 62%
		(Keum et al., 2020)	Flood prediction	Seoul, South Korea	[7.4 km ²]	м	Stormwater* - Real case	ANFIS	Rainfall(t-1), Volume (t-1), Building coverage ratio	Volume (t)	99%*	NSE (0.959)*
		(Kim et al., 2019)	Flood prediction	Gangnam area, Korea	$[7.4 \ km^2]$	м	Stormwater* - Real case	SVNARX and SOFM	Accumulative rainfall	Overflow at nodes (103)	98.50%	NSE (0.6 - 0.94)
		(She & You, 2019)	Outflow prediction	Tianjin, China	33 / [0.1314 km²]	S	Real case with synthetic data	Radial Basis function and NARX	Rainfall intensities (6)	Drainage outfall (1)	Not reported	SSE (0.273)
		(Berkhahn et al., 2019)	Flood prediction	Anonymous	1224 and 299	L, I	Stormwater* - Modifications of real cases	1 - 4 hidden layers	Precipitation intensities every 5 minutes (24 for a 2h rain event)	The maximum water level at different water cells	NA	RMSE (<0.35 cm)
		(Chiang et al., 2010)	Flood prediction	Yu-Cheng, Taiwan	[16.45 km ²]	I	Stormwater* - Real case	RNN with 1 hidden layer, 3 neurons	Registered water level and precipitation at time t (4)	Water level at time t+n (1)	NA	NSE (>0.97), CC (>0.93), NRMSE (<0.26)
	LFPB complement	(Bermúdez et al., 2018)	Surface flood volume estimation	Ghent, Belgium	6025 / [27.50 km²]	L	85% Combined - Real case	Ensemble of ANNs	Rainfall-runoff volumes aggregated over 10 and 30 min windows and volume in the underground system of the closest storage cell (3)	Presence of flooding (1) and magnitude (1)	10 ⁴ x faster*	NSE (~0.9) but variable
		(Wolfs & Willems, 2017)	Sewer water quantity simulation	Ghent, Belgium	6025 / [27.50 km ²]	L	85% Combined - Real case		Volumes between two sub- catchments (2)	Flow (1)	10 ⁶ x faster*	NSE (0.95 in average)
		(Vojinovic et al., 2003)	Wet weather flow prediction	Catchment in Auckland, New Zealand	[1.07 km ²]	S	Combined and Separated - Real case	Radial Basis function	Error, rainfall, model output (1 - 3)	Error estimation of flow (1)	NA	Improvements of 15 - 26%

193

194 *Notes:* * *denotes information not explicitly mentioned in the paper;* 'Sk' denotes a skeletonized network.

195 Acronymns: Small (S), Medium (M), Intermediate (I), Large (L); Correlation coefficient (CC). Mean squared error (MSE). Nash Sutcliff Efficiency (NSE). Root mean squared error

196 (RMSE). Mean absolute error (MAE). Squared sum of error (SSE).

197 **3.1 Metamodel Purpose**

Figure 1 shows that the two main application categories for metamodels are optimisation (48%) 198 199 and real-time applications (32%), with several examples for both WDSs and UDSs. Metamodels have been also used, although to a lesser extent, for conducting uncertainty analyses, system state 200 estimation, and to complement LFPB surrogates. The last one refers to the use of an RS method 201 (e.g., linear approximations, polynomials, ANNs) to complement an LFPB metamodel by 202 replacing a slow component or fine-tuning the outputs for better accuracy, e.g., surrogating water 203 exchange between sub-catchments with ANNs (Wolfs & Willems, 2017), or correcting the 204 205 predictions of a hydrodynamic model of wastewater flows (Vojinovic et al., 2003). In all cases, metamodels are used to reduce the computational efforts required for the hydraulic simulation of 206 207 these complex systems, which may severely compromise the feasibility of the applications.



Number of publications by application

208 209

Figure 1 Types of applications that use machine learning metamodels for Water Distribution Systems (WDS) and Urban Drainage Systems (UDS)

Optimisation usually employs population-based algorithms (e.g., genetic algorithms, particle swarm, ant colony optimisation, among others) which require multiple runs. These algorithms create an initial population, and they improve the obtained solutions through continuous iteration. Usually, these algorithms employ mechanisms inspired on genetics, such as crossover and mutation for finding (near) optimal solutions. Evolutionary algorithms are the most well-

- established metaheuristic for solving water resources problems (Maier et al., 2014); nonetheless,
- they tend to be highly computationally intensive.

Optimisation can be used to formulate and solve multiple UWN problems. This explains the high 219 number of metamodeling publications dedicated to this topic. A popular use of MLSMs for 220 optimisation in UWNs is for the (re)design of the networks. For example, applications that use 221 MLSMs include changes in pipe diameters and chlorine dosing rates (Andrade et al., 2016; Bi & 222 Dandy, 2014; Broad et al., 2005a; Sayers et al., 2019) or operation of storage tanks and pumps 223 (Broad et al., 2010; Martínez et al., 2007; Salomons et al., 2007). The goal for design is to select 224 which new system components to install or identify existing ones to substitute. For operation, the 225 aim is to find an optimal policy on how to operate the existing components. Regardless of the task, 226 the goal is to maximize the performance of the system described by the objective function(s) and 227 a number of constraints (e.g., physical, regulatory, economic, among others). In addition, other 228 problems such as water quality model calibration (Dini & Tabesh, 2017), renovation planning 229 (Dini & Tabesh, 2019), and sensor placement (Behzadian et al., 2009) have resorted to 230 metamodels. 231

Although MLSMs accelerate optimisation algorithms, they come with a series of drawbacks. First 232 of all, these models need training data (simulation examples) to calibrate their internal parameters 233 (e.g., the weights and biases in a neural network) to replicate the RS. Generating a sufficiently 234 large training dataset can be a time-consuming process, and data sufficiency depends on the 235 complexity of the input-output mapping and it can not be known a priori. Secondly, the training 236 process is another optimisation process in itself, with its own hyperparameters (e.g., learning rate, 237 number of training epochs, parameter initialization, among others depending on the optimiser) and 238 its convergence to a desired performance is not guaranteed. Furthermore, errors of approximation 239 in the RS can mislead the optimisation to suboptimal or unfeasible solutions as noted by Broad et 240 al. (2005b), especially in zones near the boundaries or outside the training range. 241

When comparing water distribution with drainage systems, it is clear that the applications of 242 optimisation in UDSs are less diverse. The reviewed papers focus on the optimisation of 243 stormwater sewers' design with Low Impact Development (LID) management (Latifi et al., 2019; 244 Raei et al., 2019; Seyedashraf et al., 2021) or for flood mitigation (Huang et al., 2015; W. Zhang 245 et al., 2019). Meanwhile, WDS optimisation is more varied, with applications to operation, 246 calibration, sensor placement, and long-term planning. This difference partially depends on the 247 stochastic nature of the rainfall events driving the functioning of combined and stormwater sewers, 248 which in turn favour real-time control over the optimisation of the operations, typical of WDS. 249 Also, the research done on MLSMs for optimisation in UDSs is rather recent (2015 or later) 250 compared to WDS (from 2005). Applications in UDSs that typically do not use metamodels can 251 benefit from the experience of tackling similar problems in the context of WDSs. Examples include 252 sensor placement (Sambito et al., 2020), calibration (Tscheikner-Gratl et al., 2016), and 253 254 optimisation of operation (van Bijnen et al., 2017).

In contrast to off-line optimisation, real-time applications require accurate answers with limited computational time. Real-time operation uses the current state of the system to modify its behaviour and improve its functioning in future time steps. In the case of UDSs, they are usually designed to retain stormwater for a certain period, to avoid combined sewer and stormwater outflows (Rosin et al., 2021; She & You, 2019) or to reduce flooding (Berkhahn et al., 2019; 260 Chiang et al., 2010; Keum et al., 2020; Kim et al., 2019; Kim & Han, 2020). Whereas, in WDSs,

the objective is to deliver high-quality drinking water while minimizing pumping costs (Pasha &

262 Lansey, 2014; Rao & Alvarruiz, 2007; Rao & Salomons, 2007).

In the case of WDSs, the reviewed real-time applications concern optimisations, in which MLSMs 263 are essential to reduce the computational time for evaluating the fitness function used by an 264 evolutionary algorithm. Consequently, these applications suffer from the drawbacks already 265 mentioned for optimisation with MLSMs. Real-time applications for UDS concern Real-Time 266 Control (RTC), where operation and validation relies on real data (Beeneken et al., 2013; 267 Langeveld et al., 2013; Lund et al., 2018). This is an issue since the usual targets are infrequent 268 events, i.e., outflows and flooding; therefore, the availability of records may be scarce or non-269 existent. 270

The third application in order of frequency is uncertainty analysis of the UWNs' performance. 271 These analyses are usually carried out via multiple simulations to test the response of the system 272 to multiple possible scenarios or uncertain input conditions, leveraging the computational 273 274 efficiency of SMs. In WDSs, ANNs have been used to replace computationally expensive models for accelerating Monte Carlo analyses. For example, Yoon et al. (2020) performed a seismic risk 275 assessment of a water distribution network considering earthquakes of different magnitudes and 276 epicentres. In UDSs, Beh et al., (2017) used metamodels to directly estimate reliability and 277 vulnerability metrics. In this case, resorting to MLSMs was crucial for the feasibility of the study. 278 Otherwise, the explicit robustness assessment would have been impossible in practice. Creating a 279 metamodel for uncertainty analysis entails having a model with explicit robustness as output, or 280 generating a training dataset with multiple runs per example. However, the former is rarely the 281 case and the latter consumes a large quantity of computational budget. 282

Other works tested the ability of ANNs to estimate the state of the system at ungauged points with 283 measurements from a limited amount of sensors. Lima et al. (2018) and Meirelles et al. (2017) 284 used recorded pressure at strategically located sensors and an ANN to estimate the pressure of all 285 the nodes in a WDS. SMs for state estimation not only decreases the degrees of freedom for the 286 addressed calibration problem but, according to the authors, they could also be used to detect 287 anomalies and predict the current state of the network in real-time. Nevertheless, in these studies, 288 289 the pressure in all the nodes is known since the MLSM is trained on simulations. For applications depending on sensor data, only a few nodes would be known and it would not be possible to 290 estimate the error for the ungauged nodes. One alternative to handle this issue is to use some 291 sensors for training and others for testing. This way, it is possible to estimate the error at the unseen 292 nodes. However, this process reduces the training data available, and it is not clear how 293 representative the testing sensors are with respect to the remaining ungauged nodes. This may lead 294 295 to unjustified trust in the model and consequent errors.

Metamodels for UDSs have also been used to complement LFPB surrogates, either to approximate some parts of the model (e.g., the most time-consuming) or to correct the predictions produced by a model. Wolfs & Willems (2017) created a modular approach in which they replaced the hydraulic simulation of drainage flow between subcatchments with an ANN, this was part of a bigger framework in which the goal was to simulate outgoing discharges for a given rainfall event. Similarly, Bermúdez et al. (2018) employed an ensemble of ANNs to accelerate a component of an LFPB model, used to estimate the occurrence and magnitude of flooding. On the other hand, Vojinovic et al. (2003) used MOUSE (MOdel for Urban Sewers), a hydrodynamic process model, to estimate flows within wastewater pipes during wet weather periods and trained a neural network to compensate for the output errors (residuals), leading to an overall increase in accuracy. Even though this hybrid approach bridges both metamodeling practices, the LPFB metamodel inherits the RS problems, e.g., database creation and training difficulties.

In summary, SMs in water networks have been primarily used for optimisation and real-time applications due to their ability to quickly evaluate outputs while remaining sufficiently accurate. This avoids running computationally expensive hydrodynamic models. Nevertheless, the use of these metamodels is not bound to these two applications. They can replace the original model for uncertainty analyses and state estimation, or help the original model by correcting outputs or approximating computationally expensive components.

314 3.2 Case studies

Figure 2 shows the number of case studies analysed in the reviewed literature. In WDSs, each 315 paper usually presents two or more networks. Since the papers introduce new problem 316 formulations or methodologies, the authors apply them to different networks to prove that the 317 methods work independently of the choice of the system. Studies in optimisation usually follow a 318 common pattern where preliminary trials are done on small benchmark networks before 319 proceeding with implementation in bigger real case scenarios. This pattern is repeated in all the 320 cases, whether it is on the same paper or in sequential papers, as in the case of the POWADIMA 321 project by Martínez et al., 2007; Rao & Alvarruiz, 2007; and Salomons et al., 2007. In the cases 322 of real-time applications, the networks were usually modified benchmarks of medium size. For 323 applications in uncertainty analysis and state estimation, the networks were real cases of large size. 324 The reviewed papers for UDSs, in contrast to WDS, present only applications with real networks, 325

some of them with modifications (e.g., Berkhahn et al., 2019; She & You, 2019).

328







On UDSs, in terms of size, most of the papers do not report the number of pipes. Consequently, 332 the extent of the system was often assessed by the reported area. This suggests that when MLSMs 333 are used, the water network is set aside and only the relation input-output is considered. The extent 334 of the case study (number of pipes or area) is a proxy of the complexity of the case studies which 335 is the relevant dimension. Nevertheless, some applications can involve medium-sized networks 336 but with high complexity (e.g., different control elements, multiple objectives, changing scenarios, 337 among others). Besides the particular characteristics of each network and application, the 338 metamodeling process was the same regardless of the size of the network. However, the required 339 time for creating the database and training the model increases with the complexity of the case 340 study. So far, the procedure does not vary as a function of the complexity of the case study; 341 342 nonetheless, considering modifications to the training process or the metamodels based on the complexity of the case study could yield better approximations to the RSs. 343

Since each system has a different area and number of pipes, we proposed the categorization in Table 1. The ratio between the number of small networks and the rest is noticeably bigger in WDSs than in UDSs due to the use of benchmarks to test the methodologies. Even though the use of metamodels is justified in larger networks, its use decreases as the size increases.

348 3.3. Metamodelling Methods

Regardless of the water network type and metamodel applications, the preferred method for 349 metamodeling is the ANN. ANNs are computational models based on the complex interaction of 350 multiple individual components (i.e., units or neurons). Each unit performs the same procedure: 351 receiving information, executing an operation (usually a linear transformation of the inputs), 352 applying a non-linear transformation to the result (e.g., hyperbolic tangent, sigmoid, rectified linear 353 unit), and sending the information to the next connected units. Each of the units has trainable 354 parameters that determine the relative weight of each of the inputs. Units are arranged in layers; 355 each ANN has at least one input layer and one output layer, where the inputs are presented to the 356 network and the computed outputs are collected, respectively. Between these layers, there are one 357 or more hidden layers, where most of the information processing takes place. ANNs learn to 358 approximate the input-output relationships in the data by tuning the trainable parameters (i.e., 359 unit's weights and biases) during the backpropagation learning process, which is usually carried 360 via gradient descent and by computing the partial derivatives of the hidden layers using the chain 361 rule of derivation. For a complete review of ANNs, the reader is redirected to Goodfellow et al. 362 (2016) for a general resource and Shen (2018) for a specific review for water resources scientists. 363

The analysis of the literature shows that the MultiLayer Perceptron (MLP) is the most widely used 364 MLSM. The MLP is a specific ANN architecture that consists of a series of layers in which all the 365 units of a layer are connected to all the neurons in the previous and next layer; hence it is also 366 known as the fully connected ANN. Most of the reviewed studies in this paper used this 367 architecture with one hidden layer; mainly due to its simplicity, high speed, and accuracy. Still, 368 the ANNs can be customized to increase the accuracy of certain applications. This practice of 369 creating deep networks, i.e., with more layers and units per layer, is part of modern deep learning 370 (Goodfellow et al., 2016). 371

In WDSs, there are two cases of variations on the number of layers: Sayers et al. (2019) used two hidden layers for optimisation of design while Yoon et al. (2020) used 15 layers in their ANN to

estimate the network performance after earthquake events. Deep networks may increase 374 performance but they are more prone to overfitting, and require more training time and examples. 375 Also, it is not possible to know the number of layers and units that yield the best performance. For 376 example, Modesto De Souza et al., (2021) tested multiple architectures of an MLP for pressure 377 estimation in a WDS. Their results suggest that the optimal number of layers is two but this can 378 vary for other applications. On the other hand, UDSs present more variation on the implemented 379 MLPs including varying the number of hidden layers (Berkhahn et al., 2019; Kim & Han, 2020; 380 Raei et al., 2019), changing the activation function to a radial basis function (She & You, 2019; 381

Vojinovic et al., 2003), and adding fuzzy logic (Keum et al., 2020).

As previously stated, MLPs are the most popular MLSM. This is not surprising due to its ease of 383 implementation and success in multiple applications, as well as hype from the AI community. 384 However, the MLP, and in general, the ML methods present several drawbacks. As Razavi et al. 385 (2012a) indicated in their numerical assessment of metamodelling strategies in computationally 386 intensive optimisation, "the likelihood that a metamodel-enabled optimizer outperforms an 387 optimizer without metamodelling is higher when a very limited computational budget is available; 388 however, this is not the case when the metamodel is a neural network. In other words, neural 389 networks are severely handicapped in limited computational budgets, as their effective training 390 typically requires a relatively large set of design sites, and thus are not recommended for use in 391 these situations.". Therefore, the use of an ANN may even harm the development of an application. 392 In that same work, the authors show that there are cases for which it is better to not use a metamodel 393 394 and go with the original model instead. Consequently, they recommend further research on determining where it is worth pursuing a metamodeling approach. In recent years, the widespread 395 availability of parallel computing (e.g., cloud computing and graphics processing unit) and user-396 friendly Deep Learning libraries, such as Pytorch (Paszke et al., 2019), have largely reduced this 397 398 problem.

Even though using MLPs is the most popular choice from the set of ML tools, it is not the only 399 one. For example, Pasha & Lansey, (2014) used support vector machines (SVMs) for improving 400 the real-time estimation of water tank levels and thus decreasing pump energy consumption in a 401 WDS. In UDSs, Chiang et al. (2010) implemented an early form of recurrent neural network 402 (RNN) for water level predictions at gauged and ungauged sites. According to the authors, their 403 decision of using this architecture was motivated by its increase in performance. However, the 404 main disadvantages of this architecture lies in training difficulty (Pascanu et al., 2013) and 405 computational costs (Strubell et al., 2020). 406

Similarly, Kim et al. (2019) and She & You (2019) leveraged the time structure in rainfall time series for real-time flood prediction with a nonlinear autoregressive network with exogenous inputs (NARX) neural networks. This architecture is a feedforward ANN that calculates the next value of a time series as a function of both past input and output values. In each study, the authors tailored the model to the conditions of their problem. Kim et al. (2019) added a second verification step to account for values that incur serious inundation damage and She & You (2019) implemented a NARX neural network for the monotonic parts of a hydrograph (i.e., ascending and descending stages) and a radial basis function MLP for the non-monotonic interval (i.e., aroundthe peak).

416 3.3.1 Metamodel inputs and outputs

The inputs to the metamodels in UWN applications are usually decision and explanatory variables 417 while the outputs can vary based on the scope of the problem. Based on the inputs used in the 418 reviewed papers, there is not a single consistent variable across the different applications in any of 419 the water networks; they are problem-specific. For example, flood prediction in UDSs relies on 420 rainfall time series, while the design of WDSs relies on inputs such as pipe diameters and chlorine 421 422 rating doses. On the other hand, the output of the metamodels are usually state variables of the UWN or performance metrics. For example, a metamodel can be developed to estimate a pressure-423 424 dependent metric, such as the resilience Network Resilience Index (NRI) (Prasad & Park, 2004), or it can output the pressures in a WDS, used to compute the NRI. Other examples of surrogated 425 components are water level in storage units or pump energy consumption. Other examples of 426 overall metrics are sampling accuracy (Behzadian et al., 2009), the economic cost of interventions, 427 428 greenhouse gases, reliability, and vulnerability (Beh et al., 2017).

Determining the output and scope of the metamodel entails deciding if the metamodel should emulate the model or one of the objectives computed after the hydraulic simulation. The reader is referred to Broad et al. (2015) for a complete methodology about metamodel scope for risk-based optimisation and its application to WDS design. In contrast, there are no applications for objective approximation using MLSMs in UDS.

By inspecting the dimensions (i.e., number) of the inputs and outputs, a converging trend is visible: 434 the number of inputs is higher than the number of outputs. This is no surprise since most of the 435 studies estimate one or two target values that summarize the desired state of the network (e.g., 436 overall performance, minimum chlorine concentration, total flooding volume) with multiple 437 decision and state variables. Nevertheless, some authors have used fewer variables to produce 438 439 more outputs. For example, in WDSs, Lima et al. (2018) and Meirelles et al. (2017) estimated 118 pressure nodes with only known pressure at 3 nodes, while Kim et al. (2019) predicted urban floods 440 441 in multiple nodes with a single rainfall time series.

On the dimensionality of ANNs, having multiple inputs and outputs allows accounting for more 442 complexity in the applications; nonetheless, they both come with downsides. For the input 443 dimensions, Razavi et al. (2012b) argue against using a large number of explanatory variables 444 (>20) since the minimum number of training examples can be excessively large. On the other side 445 of the model, the number of output variables also is recommended to be low. In theory, the number 446 of output variables is not restricted; moreover, it is one advantage of ANNs over other RS 447 metamodels as they can act as multi-output emulators. However, an ANN with multiple outputs 448 will seek to find a compromise between the errors of all the outputs, which might hurt the overall 449 accuracy of the MLSM. For this reason, an alternative approach is to train an ANN for each output 450 variable. Since each objective has a metamodel, the accuracy increases but also does the training 451 time. As noted by Andrade et al. (2016), considering one multi-output ANN or multiple ANNs 452 with single output depends on the problem at hand. The size of the water network is the most 453 important factor since, for small systems, the results with one or multiple ANNs are equivalent in 454

455 performance. In addition, the choice of one model or the other should consider desired accuracy,456 available metamodeling time, and required speed of execution.

457 3.3.2 Metamodel Performance

Regarding the performance of a metamodel, the most important characteristics are computational 458 speed and prediction accuracy. The computational saving is reported as a reduction of the time that 459 the application would have taken by running the original model. This quantity was reported by 460 nearly half of the reviewed studies and it was on average higher than 90%, most of the time over 461 98%. This is a satisfactory indication since the purpose of these SMs is to reduce the computational 462 463 burden of intensive applications. Nonetheless, around half of the studies did not report this saving. Although quantifying the computational saving is not always easy, it is recommended for future 464 researchers who use a metamodel to consider such an estimate. Since the design and training time 465 could be longer than the expected saved time, having an estimate of the potential saving aids in 466 the decision of making a metamodel. 467

In terms of prediction accuracy, there are multiple indicators used by the researchers to assess the 468 fidelity of the ML algorithm to the original model. These common metrics include root mean 469 squared error (RMSE), Nash-Sutcliffe efficiency coefficient (NSE), mean absolute error (MAE), 470 and Pearson correlation coefficient. This multitude of metrics hinders a straight comparison 471 between models or applications, but overall it is possible to observe good fittings between the 472 metamodel and the original model. It is worth noticing that the metamodel will reflect reality as 473 much as the original model is capable of doing so. Metamodels are second-level abstractions and 474 therefore may only be as good as the original model in terms of accuracy. 475

476 In addition to the previously mentioned criteria, Razavi et al. (2012b) include development time, and Asher et al. (2015) add surrogate-introduced uncertainty as assessment metrics. For these 477 criteria, seven of the reviewed papers calculated or referred to the time it took to train the models 478 and only five performed an analysis on the metamodels' robustness. Given the versatility and 479 multipurpose nature of the SMs, there are other performance indicators, e.g., ease of development, 480 explainability, generalization, or re-trainability. Along these lines, the reviewed papers disregard 481 482 these indicators since the development of the metamodel is specific for each case study and the implementation goes unnoticed. These indicators are secondary in comparison to computational 483 saving and accuracy. Both metrics constitute the most relevant metrics used in the literature, 484 including this review. 485

486 4 Current issues in metamodelling

487 Based on the current status presented in the previous section the following issues were identified.

488 **4.1. Basic applications**

MLSMs have been used to tackle various issues, namely, optimisation, uncertainty analyses, realtime applications, state forecast, and aiding LFPB metamodels. Although these generally addressed relevant problems, each of the reviewed papers had a basic framing, i.e., the inputs deal with few design or input variables (e.g., diameters, chlorine dosage, accumulated rainfall) and the outputs are usually summary variables (e.g., critical pressure, chlorine residual, flood volume). This approach is comprehensible for several reasons. First, most of the time the simplifications 495 still retain sufficient problem information to find an adequate solution. Second, it avoids problems

related to high dimensionality in the inputs and outputs. Lastly, it allows researchers to introduce
 their metamodeling method without interference from excessive complexity.

Although these frames are effective, they could result simplistic for the complexity of water 498 networks. Considering a small set of interventions may discard types and combinations of 499 interventions (e.g., allowing not only for change in diameters but also adding pumps or doing both 500 at the same time). Furthermore, other changes in the network or their components, or even 501 interactions with other city systems could be explored. However, these are rarely considered since 502 they represent a challenge for traditional RS metamodels; current MLSMs are very specific to the 503 cases in which they are trained on. Because of this, new approaches are required, mainly in 504 optimisation and uncertainty analysis. 505

As seen in section 3, the most popular application for MLSMs is optimisation. In this application, multiple authors (Beh et al., 2017; Doorn, 2021; Kapelan et al., 2005; Razavi et al., 2021) have remarked on the importance of considering new objectives. For example, robustness for designing water systems, especially under deep uncertainty, requires considering multiple scenarios for which is not possible to assign a probability or ranking. This analysis is desirable because water networks are systems with long lifespans of service. Nonetheless, objectives like robustness tend to be more computationally intensive; therefore, their need for metamodels increases.

A relevant missing layer of complexity is uncertainty analysis, especially for UDSs. The current 513 practice to design the system is to use a single benchmark storm and assume it is representative of 514 the future rain events the system will face. However, two UDSs with similar performance during 515 a design event could behave very differently for other rainfall patterns. According to Ng et al. 516 (2020), the final design considering a single strong storm does not guarantee optimal performance 517 during long mild storms and for a succession of frequent small events. Naturally, the authors 518 519 recognize that performing a design considering multiple events would increase the computational effort but also suggest the implementation of SMs for dealing with this difficulty. 520

521 **4.2 Case studies: Lack of benchmarking with complex networks**

522 Benchmark water networks are open access datasets that contain the necessary information to 523 create models of a system. It consists of the topology of the network, its components, and 524 depending on the system it could incorporate leakages, demand patterns, cyber-attacks, rainfall, or 525 surveillance data. Benchmarks are used as reference points to compare the performance of models 526 and algorithms. Here, it is necessary to distinguish between synthetic and real data. Even though 527 the synthetic data allow to implement and compare algorithms, they may not reflect all the 528 processes that real data can account for.

There is a clear difference between types of infrastructure in the number of used networks since benchmark networks in UDSs are not as available as in WDSs. In water distribution, there is a set of water networks called Water Distribution System Research database. The ASCE Task Committee on Research Databases for WDS created this database which is hosted by the University of Kentucky (2013). There are benchmarks for multiple problems in categories such as network expansion, operation, and design. This allows modellers to easily obtain data for the development and comparison of algorithms in networks of different sizes. On the other hand, there 536 is no consolidated set of benchmark networks for UDSs, let alone an entire structured database. 537 This is attributable to factors such as the difficulty of taking measurements in sewer environments 538 and, according to Pedersen et al. (2021), the little interest of utility companies in making the 539 datasets publicly available. Consequently, all the applications on UDSs were entirely developed 540 for real cases, which is positive for the bridging between the theoretical approaches and the 551 practice, but hampers the development of algorithms on the systems, due to the difficulty of 552 comparison and the process of accounting for particularities of each system.

Regarding the size of the case studies, most of the systems in which the MLSMs were used were 543 medium or small. Metamodels are most useful in problems with large computational times, that 544 is, in applications with large water networks. In the case of WDSs, a common practice to test the 545 effectiveness of a method is developing a metamodel for a small benchmark network and then 546 using the same steps for creating a metamodel in a big real case. Even though this practice is 547 reasonable, it assumes the response surface of both networks is comparable or similar. However, 548 this is not necessarily the case as reported by Andrade et al. (2016) who noted contrasting 549 accuracies between big and small case studies when training metamodels. Exploring solution 550 spaces is already an issue when using metamodels, independent of the network, as reported by 551 Broad et al. (2005), but large networks represent additional challenges that increase in complexity 552 in a non-linear manner. 553

4.3 Machine learning and multi-layer perceptron limitations

Although the MLP is not the only ML technique, it is the most popular one among MLSMs. Given that its structure allows it to address multiple types of problems, it has become a one-size-fits-all model. Nevertheless, it presents multiple issues, namely, the curse of dimensionality, black-box nature, and rigid structure. These three shortcoming respectively 1) hinder their use for high dimensionality problems, 2) limit confidence in their approximations, and 3) prevent the transferability of trained models across different case studies.

4.3.1 Curse of dimensionality - Metamodeling time

The curse of dimensionality indicates that for a certain level of accuracy, there is an exponential 562 increase in the required amount of data as the dimensions of a problem increase (Keogh & Mueen, 563 2017). Naturally, this problem can be addressed by reducing the number of input dimensions (i.e., 564 fewer explanatory variables) using prioritization based on experience, knowledge of the task, or 565 some automatic procedure such as principal component analysis (PCA). However, as noted by 566 Maier et al. (2014), for real-world problems reducing the number of input features may not be a 567 satisfactory solution because it usually leads to an approximation that could exclude optimal zones 568 and prevent the algorithms to find optimal solutions. Given this situation, searching for solutions 569 on the algorithmic side may yield better answers. 570

The SMs have worked adequately so far but future metamodels are likely to increase in complexity.
This is either due to an increase in the complexity of UWNs or an increase in the number of input
(more design choices/explanatory variables) or output (more objectives) dimensions. Both drivers

574 increase the size of the metamodels and consequently the number of training examples. Since the

original models are already expensive to run, creating a large training dataset might be unfeasible

576 in the first place. The metamodeling time would become the obstacle. This time is usually

577 disregarded since some authors consider it not relevant compared to the posterior computational

578 gain in the application. Nevertheless, this time is important in high dimensional search spaces, as 579 noted by Razavi et al. (2012b), since the number of design samples required to train the metamodel

579 noted by Razavi et al. (2012b), since the number of design samples re
580 could be already prohibitively large.

581 4.3.2 Black box nature - Deterministic and obscure outputs

Two of the most recurrent criticisms of ML models are their lack of uncertainty estimation and the lack of their transparency, i.e. little or no ability to explain the results they obtain. Both are overlooked aspects of metamodeling in the context of UWNs. The MLSMs return a unique answer without uncertainty bands or possibilities to explain the combination of inputs that drove to the final outputs. For SMs, these issues are not major concerns; nevertheless, their inclusion aids the applications in which the SMs are used.

Regarding uncertainty estimation, a few papers (Raei et al., 2019; Rosin et al., 2021; She & You, 588 2019; W. Zhang et al., 2019) estimated the effect of including a metamodel in their respective 589 application. Not accounting for this uncertainty can lead to bad approximations of the actual 590 response surface and suboptimal or unfeasible solutions. Authors have dealt with this difficulty by 591 performing sensitivity analysis (e.g., Raei et al., 2019) or training multiple models in parallel with 592 slightly different datasets and averaging the outputs of the models. For example, Rosin et al. (2021) 593 developed a committee of ANNs with this approach. However, this analysis requires extra 594 considerations which may increase the metamodeling time. Some guidelines have been given for 595 the pre-treatment (Broad et al., 2015) and post-treatment (Broad et al., 2005a) of these SMs but 596 there is still a lack of focus on improving the management of uncertainty during treatment, i.e., 597 developing a model that directly considers uncertainty. Algorithms in the branch of robust ML 598 may contribute to aid in the direct incorporation of metamodel uncertainty quantification whether 599 it comes from the data (Wong & Kolter, 2019) or the model (Loquercio et al., 2020). 600

Although robust learning allows estimating the uncertainty of a result, it cannot explain why. This 601 is the area of explainable ML. For water networks' SMs, being able to explain the results would 602 help to understand the relationship between the decision variables and the objective function for 603 the particular network that is being surrogated. For example, understanding which pipes (or a 604 combination of them) play a key role in the resilience or flooding in a water network. There is a 605 growing interest in the AI community towards explainable models to gain insights (Bhatt et al., 606 2020), ensure scientific value (Roscher et al., 2020), and develop trust in the outcomes of ML 607 608 models (Dosilovic et al., 2018).

609 4.3.3. Rigid architecture - Specific case use

One disadvantage of MLSMs is the high degree of specialization in the trained metamodel. As

seen before, these metamodels achieve high accuracies in the data for which they were trained.

However, once they are trained, they become specific and rigid. Their structure limits its use for

other tasks in the same system or similar applications in other water networks. The metamodel can be run several times on the same water network but doing the same operation in a different system requires a new metamodel, which should be trained from scratch. This is not desirable since the training process could consume most of the computational budget, especially in large case studies.

One solution is to leverage the training process of other models with transfer learning to decrease 617 the number of examples to train a new model. Situations for which transfer learning is desirable 618 are changes in the water network composition, similar system metamodeling, and change in the 619 behaviour of the surrogated system. Changing components of the system accounts for scenarios 620 when components (e.g., pipes, pumps, or tanks) are added to or removed from the system. Even 621 though the system changes, it is still related enough to leverage a pre-trained model on that water 622 network. In a similar way, two networks can share enough resemblance (e.g., a subsystem of 623 another network, two skeletonized networks, or two networks with similar topology or geography) 624 that it makes sense to use an SM from one as a pre-trained SM for the other. Lastly, when the 625 system changes and the metamodel no longer applies is a challenge, also known as concept drift, 626 that can be addressed using transfer learning. Here the two related water networks are the same 627 but in two different periods. 628

629 4.4. Gaps in Knowledge

Based on the above critical analyses of metamodels and the issues identified the following key gaps in knowledge are summarised here:

Lack of depth on optimisation of complex objectives and uncertainty analysis for water
 networks using MLSMs. There are still additional and more complex objectives that can be
 optimised with the aid of MLSMs, for instance, robustness and interventions under deep
 uncertainty.

Lack of benchmark water networks, especially for UDSs and complex cases. First, this
 hinders the development and comparison of algorithms across studies, and second, these
 metamodels still lack research on the changes of the response surface with the increase in the
 complexity of the water system, especially for large systems

Gurrent MLSMs' limitations prevent advanced metamodeling applications. MLSMs can
easily grow in size when the complexity of the response surface increases, most of the applications
do not consider the uncertainty added by the metamodel, and its structure makes it rigid and not
(re)usable for other cases.

6445 Research directions

Based on the identified gaps, three main lines for future research are suggested. They consider the current and future needs in applications on UWNs as well as the potential of MLSMs to meet them.

647 5.1 Advanced applications

The current needs for adaptable water infrastructure are based on drivers such as growing demographics, urbanization, and climate change. As indicated in the UN-Water report "Water and Climate Change", taking adaptation and mitigation measures benefits water resources management and improves the provision of water supply and sanitation services. In addition, it contributes to combat both causes and impacts of climate change while contributing to meeting 653 several of the Sustainable Development Goals (UNESCO, 2020). In UWNs, multi-objective 654 optimisation and uncertainty analysis play a key role in the search for adaptation measures and 655 decision making, and MLSMs can help improve and accelerate their implementation.

Optimisation applications will increase in the number and complexity of the inputs and outputs. 656 Increasing the number of inputs, i.e., decision variables and design interventions (e.g., nature-657 based solutions), allows to explore more alternatives, consider uncertainty, or assess multiple 658 scenarios. On the other hand, the output of the optimisation is leaning towards complex objectives 659 such as multi-objective robustness (e.g., Kasprzyk et al., 2013), multiple technical performance 660 metrics (e.g., Fu et al., 2013), pro-active maintenance (Kumar et al., 2018), complex water quality 661 indicators (Jia et al., 2021), and human values (Doorn, 2021). Multi-objective optimisation allows 662 identifying solutions balancing trade-offs among objectives, for instance, cost and resilience 663 (Wang et al., 2015). Naturally, when considering more objectives, the computational load 664 increases, especially when those objectives are computationally expensive (e.g., robustness). In 665 previous phases of research on optimisation, metamodels were seen as an aid, but as optimisation 666 gradually evolves to consider additional and more complex objectives, metamodels become 667 indispensable (e.g., Beh et al., 2017). 668

Regarding uncertainty analysis, it is necessary to have fast, reliable, and flexible metamodels that 669 can adapt to the multiple conditions in which the systems are evaluated and under multiple criteria. 670 Traditionally, simplified models have been preferred for this task; however, RS metamodels 671 become appealing alternatives when dealing with more complex objective functions and original 672 models. Metamodels should play a key role in the development of frameworks for robustness-673 driven design. This application has major implications for UDSs, since no MLSM study focused 674 on uncertainty analysis, even when the evidence suggests the criteria for the design of these 675 systems is not necessarily robust (Ng et al., 2020). Although uncertainty analysis entails an 676 intrinsic increase in the computational effort, the benefits they bring outweigh the challenges it 677 represents. According to the IPCC (2021b), UDSs are expected to receive more intense rainfall 678 events based on climatic projections but considerable uncertainty remains. 679

The community should further research combined RS-LPFB applications, to further integrate MLSMs with physically-based models for accelerating the underlying hydrodynamic engines. Likewise, physically-based models could be hybridized by incorporating an ML model that corrects the outputs of the original model for higher accuracy accounting for the real behaviour of the system. Looking ahead, ML algorithms could detach from the physically-based model and replace its functioning with a cheaper version to run based on increasingly available real-world data (e.g., digital twins for UWNs (IWA, 2021)).

5.2 Benchmarking and large network behaviour

The lack of benchmark models is a gap that was already identified by Maier et al. (2014) who set the characteristics and recommendations of valuable benchmarks, including non-trivial real-world problems with a representative range of decision problems characteristic of the water systems. The review shows that UDSs lack such benchmarks. To overcome this issue, we recommended to implement a similar approach to that of the Kentucky database, with applications such as real-time control, outflow, and flood prediction. For WDSs, it is appropriate to enlarge the current databases to account for new objectives, interventions, performance metrics, and real case examples. Regarding metamodels, the benchmarks should also include a reference model to compare computational saving and accuracy, with suggested performance metrics, such as NSE, RMSE, or the number of model executions.

As Goodfellow et al. (2016) indicate, having benchmark databases with real cases is one of the 698 reasons why deep learning has recently become a crucial technology in several disciplines. In AI, 699 datasets went from hundreds or thousands of examples in the early 1980s up to datasets with 700 millions of examples after 2010. Nowadays, thanks to the increase in connectivity and 701 digitalization of our society, a large amount of ML algorithms can be fed with the information they 702 require to achieve high accuracy. Since the ML and DL models are dependent on their training 703 sets, their success goes hand in hand with the size and quality of available datasets, preferable with 704 real information. The UWNs' research community is moving the first steps in this direction. One 705 example concerns the UDS of the Bellinge dataset (Pedersen et al., 2021), a suburb to the city of 706 Odense, Denmark that is now available for "independent testing and replication of results from 707 future scientific developments and innovation within urban hydrology and urban drainage system 708 research". This dataset includes 10 years of asset data (information from manholes and links), 709 sensor data (level, flow, and power meters), rain data, hydrodynamic models (MIKE urban and 710 EPA SWMM), and other information. Similar examples are needed to enable the exploration of 711 metamodels' responses in networks of different characteristics (e.g., size, connectivity, slope). 712

713 As for the size of the networks, further research is required to assess the response surface of large networks. Specifically, new benchmark datasets should also include complex network cases for 714 715 their study. These can be large networks or medium-size cases with high complexity. Considering that the larger the network the higher the required time to generate and use the training data, 716 significant efforts are required on this matter. Metamodels could aid in reducing the computational 717 times that obstruct studying the response surface of large and complex systems. Nonetheless, new 718 metamodels are required to account for the complexity of these cases and use as few training 719 scenarios as possible. 720

5.3 Unexplored advanced metamodeling technologies

ML is the area with the highest growth in academic output in recent years. However, the field of MLSMs for UWNs has not yet considered the new tools and algorithms recently developed by researchers in fundamental AI or other applied disciplines. These advancements include DL architectures that express assumptions of the data in the ANNs for robust, interpretable, and transferrable models. This new wave of AI formalizes the attempts to add knowledge about modelled processes as well as extract knowledge from the results.

5.3.1 Inductive bias – Deep learning: Graph Neural Networks

The curse of dimensionality can be addressed by including inductive biases. Following the work of Battaglia et al. (2018), we define the inductive bias as the "expression of assumptions about either the data-generating process or the space of solutions". Inductive bias can be seen as well in the architecture of the model by leveraging the inner structure of the data, which could be spatial, temporal, or relational. Exploiting the structural information of the data can reduce the number of parameters, and consequently the required training examples by parameter sharing and sparsity of connections. The data structure gives information about the similarity of the data points in a relevant dimension (e.g., distance, time, connection). In that sense, similar data can be treated analogously (parameter sharing) and dissimilar data can remain unrelated (sparse connectivity).

Inductive bias nudges a learning algorithm to prioritize some solutions over others. This allows 738 finding high-performing solutions more easily than when it is not considered. Ideally, involving 739 inductive bias improves the search for solutions without compromising the performance, as long 740 as the right inductive bias is chosen; otherwise, it can lead to suboptimal performance (Battaglia 741 et al., 2018). For example, when surrogating the pressure at the nodes of a WDS with a neural 742 network (e.g., Broad et al., 2005; Meirelles et al., 2017) there are multiple metamodel solutions, 743 i.e., architectures with specific parameter values that can approximate the response surface 744 described by the training data. Nevertheless, when adding inductive bias, the set of possible 745 solutions shrinks to a subset of solutions that comply with predefined characteristics, for example, 746 having graph structure, following physical laws, or agreeing with measurements. 747

The most common components in DL are fully connected, convolutional, recurrent, and, more 748 recently, graph layers. The fully connected layers have a weak inductive bias, while each of the 749 remaining exploits some relation or invariance in the data. The convolutional layers typical of 750 convolutional neural networks (CNNs) leverage the regular structures in grids, such as images, 751 and connects information according to Euclidean closeness. Recurrent neural networks (RNNs) 752 consist of recurrent units which consecutively process data sequences, such as time series, and 753 connects information according to sequential similarity. On the other hand, graph neural networks 754 (GNNs) extend DL methods to non-Euclidean data, such as graphs, where entities are connected 755 756 by relations or, in graph terminology, nodes connected by edges.

Given their relational inductive bias, GNNs are the most suitable DL architecture for applications 757 in UWNs, since the natural structure of these systems is a graph. Researchers have already 758 exploited graph theoretical concepts to develop decomposition models of WDNs (Deuerlein, 759 760 2008), assess the resilience of sectorized WDNs (Herrera et al., 2016), as well as identifying critical elements in UWNs (Meijer et al., 2018, 2020). Furthermore, there are already some 761 applications of GNNs in UWNs. In WDSs, Tsiami & Makropoulos, (2021) employed this 762 architecture for cyber-physical attack detection using a graph created from sensors in the water 763 system. In UDSs, Belghaddar et al. (2021) applied this method to database completion of 764 wastewater networks. 765

This architecture operates on the graph domain, which allows it to leverage the pre-existing 766 767 network topology of the data. This architecture has gained considerable attention in the last years due to its ability to include relational structure from connected entities. Even though GNNs' 768 outputs continue to be hardly explainable, there are efforts to generate explanations of their 769 outputs, e.g., GNNExplainer (Ying et al., 2019). As noted by Battaglia et al., (2018), "the entities 770 and relations that GNNs operate over often correspond to things that humans understand (such as 771 physical objects), thus supporting more interpretable analysis and visualization". In this way, 772 773 GNNs are not entirely explainable but they are more explainable than other DL architectures.

It is also possible to use combinations of layers in problems that contain more than one structure such as in the case of UWNs, which have temporal, spatial, and topological variability. An example

of the application of these graph models in a civil infrastructure was developed by Sun et al. (2020) who included the spatial and temporal relations in a road network for traffic forecasting. This infrastructure has multiple parallels with UWNs, including its graph connectivity, spatial-temporal
 variability, and human interaction. Another similar infrastructure with more examples can be
 found in power systems for which GNNs have been used in key applications such as fault scenario

application, time series prediction, power flow calculation, and data generation (Liao et al., 2021).

For a review in depth of GNN architecture, the reader is referred to Zhou et al. (2018).

This architecture presents an opportunity to leverage the present structure of the data generated in 783 the UWNs to decrease the number of parameters and consequently the required training data; 784 which enables creating SMs of larger networks and many and more complex objectives. By 785 conditioning the characteristics of the solutions, the metamodels gain the possibility to generalize 786 to similar cases. For example, pipe changes in a network configuration could be better represented 787 with a GNN-based metamodel. This GNN SM could be able to adjust itself without modifying the 788 underlying structure, which would probably be required in the case of other metamodels that do 789 not consider this inductive bias. 790

791 5.3.2 Third wave of Artificial Intelligence

The US Defense Advanced Research Projects Agency (DARPA, 2016) separates the different phases of AI into three waves. The first wave refers to the past approaches and the birth of AI, the

second wave is the current and popular phase of high-performing black boxes, and lastly, the third

795 wave is proposed for the future of AI with models leaning towards robustness and explainability.

Robustness refers to the ability to include uncertainty in the calculation of the outputs of a model, 796 in this way the user not only receives a deterministic answer but a range of possible values, usually 797 represented by an expected value (e.g., mean) and a measure of uncertainty (e.g., variance). 798 According to Gawlikowski et al. (2021), methods for estimating uncertainty in ANNs can be split 799 into four types: single deterministic methods, bayesian methods, ensemble methods, and test-time 800 augmentation methods. Each of these lines offers an estimation of the degree to which the neural 801 network is certain of the output. This aspect is relevant when quantifying how likely it is for the 802 803 metamodel to detach from the response surface which may cause, depending on the application, to omit optimal solutions, miss outflows, or underestimate floods. Recommended methods for 804 implementation on MLSMs include Bayesian neural networks (e.g., Zhu & Zabaras, 2018) or 805 single deterministic methods, the latter is recommended based on the low additional computational 806 burden they include. 807

Research in explainability has also gained popularity in recent years. In the case of MLSMs, having 808 an explainable model would allow us to better understand the response surface of the original 809 model or the solution space. An improved comprehension of the response surface would facilitate 810 obtaining a better insight on the behaviour of different algorithms (e.g., evolutionary methods); 811 ultimately, contributing to what type of heuristic is best suitable in each application in water 812 network which is a topic in which we have still very little understanding of (Maier et al., 2014). 813 On the other hand, solution space explanation would allow gaining insight about which and 814 components in the real system affect its performance, but most importantly, how they affect it. 815 816 This could drive the interventions in the physical water network to improve its performance. Recommended models for implementation in this category are GNNs, as already reported by 817 Tsiami & Makropoulos (2021), who were able to perform a removal analysis to quantify the 818 819 contribution of each considered component (e.g., valves, tanks, and pumps) of the physical water

- network to the model's performance. Since GNNs' structure resemble the underlying system, it is
- 821 possible to relate events on the metamodel to the actual system.
- 822 5.3.3 Transferrable AI models

The reviewed studies in this paper presented a methodology for training a metamodel to surrogate a computationally expensive model. Although the methodology is transferrable, meaning the steps can be followed and repeated to obtain a similar metamodel in another case study, the metamodel itself cannot be transferred to a new case study. This implies that all the metamodeling time spent on training is specific for every case. Through transferrable models, the authors may develop not only methodologies but also pre-trained SMs, which can be adapted to other cases lowering the amount of training needed for this new network.

Having a transferrable model would allow training the metamodel with data not only from the case 830 study at hand but also from other, real and synthetic cases. For example, the benchmark datasets 831 discussed previously. This increase in available information to train on is expected to improve the 832 performance of the metamodel or even allow it to exist for cases in which data is scarce, for 833 example, very computationally expensive UWNs in which training examples are costly. Once 834 again, inductive bias plays a role, since the assumptions added to the algorithm delimit a smaller 835 solution space, the ML models can be used as pre-trained solutions for other tasks. In the AI 836 domain, this practice is referred to as transfer learning. Transfer learning is mainly implemented 837 for specialized deep learning methods, i.e., architectures with strong inductive bias. It has been 838 successfully implemented for applications such as diagnosis of medical images using CNNs 839 (Vogado et al., 2018), prediction of air pollutants using RNNs (Hang et al., 2020), and 840 bioinformatics as well as social-network classification tasks with GNNs (Verma & Zhang, 2019), 841 among others (Weiss et al., 2016). 842

For transferrable SMs in UWNs, GNNs seem to be the natural option based on the agreement 843 between the structure of the real system and the inductive bias corresponding to the GNNs. In an 844 analogous way that CNNs learn filters that are independent of the input (i.e., images), GNNs learn 845 filters that can be used across cases (e.g., water networks). Adding the structure and physics to the 846 847 metamodel allows including more domain knowledge in the ANN that improves generalization capabilities. A relevant example of a model like this is the mass conserving RNN for rainfall-848 runoff modelling developed by Hoedt et al. (2021) in which the parameters used in the model 849 resemble the mass conservation principle, which increased the accuracy and improved the model's 850 interpretability. At the same time, transferability opens the door to new applications, such as online 851 optimisation of interventions, by learning the effect of changes in the topology and components of 852 the network. 853

Using physical information, such as the knowledge embedded in the hydrodynamic models, also 854 allows generating hybrid and general models. These models allow bridging the best of two 855 domains: physical-based and data-driven. On this, Vojinovic et al. (2003) indicated that "the major 856 advantage of integrating both a deterministic (numerical) model and a stochastic (data-driven) 857 858 model over using the stochastic data-driven model alone is that the already available deterministic model quality is exploited and improved, instead of starting from scratch and throwing away all 859 knowledge." Furthermore, combining the domain knowledge with transferable models opens the 860 possibility of creating general models. This type of model detaches from the training set in which 861

it was trained so that its predictions can be applied in unseen scenarios. Following this trend, 862 Kratzert et al. (2019) developed a recurrent ANN trained on basins from a continental dataset using 863 meteorological time series data and static catchment attributes, and they were able to outperform 864 hydrological benchmark models calibrated on individual catchments. The analogous application 865 in UWNs would be an ML-based hydrodynamic model trained on a set of distribution or drainage 866 systems which can generalize to independent unknown water networks. Such "DeEPANET" or 867 "DeepSWMM" models can be developed by leveraging the inductive bias of GNNs, and 868 accounting for the time dimension with recurrent layers or by resorting to an encoder-decoder 869 architecture (Du et al., 2020). 870

871 6 Conclusions

This work reviews the current state of the application of MLSMs in urban water networks and proposes promising forward directions based on recent and successful developments in ML.

In terms of purpose, the main uses of MLSM in UWNs are optimisation and real-time problems. 874 Even though MLSM accelerate optimisation algorithms by increasing the speed of individual 875 iterations, these algorithms have multiple disadvantages. The training process can be time-876 consuming and the required size of that dataset cannot be known a priori as it depends on the 877 complexity of the input-output mapping. For case study type, the UWNs in which MLSMs are 878 applied vary in size and type. For analysing the complexity of the case studies, we preffered to 879 consider WDSs and UDSs separately. Regarding its use in WDSs, the papers follow a clear pattern: 880 the development and trial are usually made in medium or small benchmark networks, and the 881 posterior implementation of the metamodel is done in a large real network. On the other hand, 882 UDSs do not count with applications on benchmark networks due to their lack of availability. In 883 terms of the metamodel, except for some applications of SVMs or RNNs, the vast majority of 884 applications used MLP as SM. This method has been successfully implemented due to its high 885 accuracy and flexibility regarding the inputs and outputs that it can map. Nevertheless, the MLSMs 886 present multiple drawbacks that may even harm the development of an application. It is advisable 887 to consider if an MLSM is worthwhile before starting its training. 888

Based on the reviewed literature, the following issues and gaps in knowledge were identified in
terms of limitations of existing MLSMs. These problems include limitations on the MLSMs, lack
of depth in current applications, and insufficient benchmarking datasets.

- Regarding metamodels' limitations, current MLSMs have the following issues: they can easily grow in size when the complexity of the response surface increases, most of the applications do not consider the uncertainty added by the metamodel, and its structure makes it rigid and not (re)usable for other cases.
- In terms of applications, optimisation is where most of the SMs are currently used;
 nevertheless, there are still additional and more complex objectives that can be
 optimised with the aid of MLSMs, for instance, robustness and interventions under
 deep uncertainty.
- On case studies, the reviewed papers denote two main issues: first, there is a lack of UDSs benchmarks, which hinders the development and comparison of algorithms

- 902across studies, and second, these metamodels still lack research on the changes of the903response surface with the increase in the complexity of the water system, especially for904large systems.
- The following research directions are suggested to address the above key gaps in knowledge:
- 906 • Regarding metamodeling methods, further research is required on advanced 907 metamodeling techniques that include: inductive bias, robustness, and transferability. The notion of inductive bias allows leveraging prior information to reduce the required 908 training samples. Examples of this bias include adding physical laws, coherence with 909 910 sensor data, or considering the underlying structure of the data - space, time, or topology– In this regard, the recently developed GNNs resemble the already existing 911 912 architecture of the urban water networks and offer the highest fit to the data in these systems. Furthermore, the new approach for AI models is to focus on the robustness 913 and explainability of the models which offer insight into the applications and 914 opportunities for improvement in the actual systems. Moreover, implementing the new 915 916 architectures of ML as an SM would allow transfer learning, which represents the ability to use pre-trained models and save computational budget. 917
- On applications, additional efforts are encouraged in two areas in which metamodels will increasingly be more required: uncertainty analysis and multi-objective optimisation, especially when robustness metrics are used as optimisation objectives.
 Further research is required on other less developed applications, namely, real-time predictions, state estimation, and to a lesser extent, LFPB complements. These applications have been minimally explored and most of them have only been used for a specific type of water network.
- Regarding case study type, it is crucial to develop benchmark UWNs, especially of UDSs, and complex networks. This data will facilitate training, testing, and comparing new metamodels. These new benchmarks could incorporate information on leakages, demand patterns, cyber-attacks, rainfall, or surveillance data as well as performance metrics as reference points to compare performance.

Exploring the potential of MLSMs for approximating UWNs' components and correcting predictions with real data can lead to independent ML models of the water networks that leverage the physical domain knowledge and the measurements. New MLSMs are encouraged to leverage the inductive bias offered by the increasing data to help UDS and WDS operators. The new advancements in ML, especially GNNs, have great potential to advance surrogate modelling in UWNs. Water network modellers can speed up calculations for larger and more complex cases, being able to design more robust and overall better urban water systems.

938 List of abbreviations and acronyms

- 939 AI Artificial intelligence
- 940 ANN Artificial neural network
- 941 CNN Convolutional neural network
- 942 DL Deep learning
- 943 GNN Graph neural network
- 944 LFPB Lower-fidelity Physically-based
- 945 MAE Mean absolute error
- 946 ML Machine learning
- 947 MLSM Machine learning-based surrogate model
- 948 MLP Multi-layer perceptron
- 949 NSE Nash-Sutcliffe efficiency coefficient
- 950 RS Response surface
- 951 RMSE Root mean squared error
- 952 RNN Recurrent neural network
- 953 SM Surrogate model
- 954 SUDS Sustainable urban drainage systems
- 955 UDS Urban drainage system
- 956 WDS Water distribution system.
- 957 Acknowledgments
- This work is supported by the TU Delft AI Labs programme.
- 959 **Open Research**
- No experimental data or code were produced for this manuscript.

961

963 Authors Contribution

- All authors contributed in conceptualising the review paper and its outline. AG wrote most of the
- paper, produced all figures and tables, and formatted the article. RT wrote parts of the paper across
- all sections. ZK proposed the initial idea. RT, ZK, and JL reviewed, revised, and supervised the
- 967 progress of the paper.
- 968 **References**
- Andrade, M. A., Choi, C. Y., Lansey, K., & Jung, D. (2016). Enhanced artificial neural networks
- 970 estimating water quality constraints for the optimal water distribution systems design.
- 971 Journal of Water Resources Planning and Management, 142(9).
- 972 https://doi.org/10.1061/(ASCE)WR.1943-5452.0000663
- Asher, M. ., Croke, B. F. ., Jakeman, A. ., & Peeters, L. J. . (2015). A review of surrogate models
 and their application to groundwater modeling. *Water Resources Research*.
- 975 https://doi.org/10.1029/eo064i046p00929-04
- 976 Battaglia, P. W., Hamrick, J. B., Bapst, V., Sanchez-Gonzalez, A., Zambaldi, V., Malinowski,
- 977 M., Tacchetti, A., Raposo, D., Santoro, A., Faulkner, R., Gulcehre, C., Song, F., Ballard, A.,
- Gilmer, J., Dahl, G., Vaswani, A., Allen, K., Nash, C., Langston, V., ... Pascanu, R. (2018).
- Relational inductive biases, deep learning, and graph networks. *ArXiv*, 1–40.
- Baú, D. A., & Mayer, A. S. (2006). Stochastic management of pump-and-treat strategies using
- surrogate functions. *Advances in Water Resources*, 29(12), 1901–1917.
- 982 https://doi.org/10.1016/J.ADVWATRES.2006.01.008
- Beeneken, T., Erbe, V., Messmer, A., Reder, C., Rohlfing, R., Scheer, M., Schumacher, B.,
- Weilandt, M., Weyand, M., Erbe, V., Messmer, A., Reder, C., Rohlfing, R., & Scheer, M.
- 985 (2013). Real time control (RTC) of urban drainage systems A discussion of the

- 986 *additional efforts compared to conventionally operated systems.* 9006.
- 987 https://doi.org/10.1080/1573062X.2013.790980
- Beh, E. H. Y., Zheng, F., Dandy, G. C., Maier, H. R., & Kapelan, Z. (2017). Robust optimization
- 989 of water infrastructure planning under deep uncertainty using metamodels. *Environmental*
- 990 *Modelling and Software*, 93, 92–105. https://doi.org/10.1016/j.envsoft.2017.03.013
- Behzadian, K., Kapelan, Z., Savic, D., & Ardeshir, A. (2009). Stochastic sampling design using a
- 992 multi-objective genetic algorithm and adaptive neural networks. *Environmental Modelling*

993 *and Software*, 24(4), 530–541. https://doi.org/10.1016/j.envsoft.2008.09.013

- Belghaddar, Y., Chahinian, N., Seriai, A., Begdouri, A., Abdou, R., & Delenne, C. (2021). Graph
- convolutional networks: Application to database completion of wastewater networks. *Water* (*Switzerland*), 13(12), 1–19. https://doi.org/10.3390/w13121681
- 997 Berkhahn, S., Fuchs, L., & Neuweiler, I. (2019). An ensemble neural network model for real-
- time prediction of urban floods. *Journal of Hydrology*, 575, 743–754.
- 999 https://doi.org/10.1016/j.jhydrol.2019.05.066
- 1000 Bermúdez, M., Ntegeka, V., Wolfs, V., & Willems, P. (2018). Development and Comparison of
- 1001 Two Fast Surrogate Models for Urban Pluvial Flood Simulations. *Water Resources*
- 1002 *Management*, 32(8), 2801–2815. https://doi.org/10.1007/s11269-018-1959-8
- 1003 Bhatt, U., Xiang, A., Sharma, S., Weller, A., Taly, A., Jia, Y., Ghosh, J., Puri, R., Moura, J. M.
- 1004 F., & Eckersley, P. (2020). *Explainable Machine Learning in Deployment*. 648–657.
- 1005 Bi, W., & Dandy, G. C. (2014). Optimization of water distribution systems using online retrained
- 1006 metamodels. Journal of Water Resources Planning and Management, 140(11).
- 1007 https://doi.org/10.1061/(ASCE)WR.1943-5452.0000419
- 1008 Broad, D. R., Dandy, G. C., & Maier, H. R. (2005a). Water distribution system optimization

- 1009 using metamodels. Journal of Water Resources Planning and Management, 131(3), 172–
- 1010 180. https://doi.org/10.1061/(ASCE)0733-9496(2005)131:3(172)
- 1011 Broad, D. R., Dandy, G. C., & Maier, H. R. (2005b). Water Distribution System Optimization
- 1012 Using Metamodels. Journal of Water Resources Planning and Management, 131(3), 172–
- 1013 180. https://doi.org/10.1061/(asce)0733-9496(2005)131:3(172)
- 1014 Broad, D. R., Dandy, G. C., & Maier, H. R. (2015). A systematic approach to determining
- 1015 metamodel scope for risk-based optimization and its application to water distribution
- 1016 system design. *Environmental Modelling and Software*, 69, 382–395.
- 1017 https://doi.org/10.1016/j.envsoft.2014.11.015
- 1018 Broad, D. R., Maier, H. R., & Dandy, G. C. (2010). Optimal Operation of Complex Water
- 1019Distribution Systems Using Metamodels. Journal of Water Resources Planning and
- 1020 *Management*, *136*(4), 433–443. https://doi.org/10.1061/(asce)wr.1943-5452.0000052
- 1021 Brown, R. R., Keath, N., & Wong, T. H. F. (2009). Urban water management in cities: historical,
- 1022 current and future regimes. *Water Science and Technology*, 59(5), 847–855.
- 1023 https://doi.org/10.2166/wst.2009.029
- 1024 Chiang, Y.-M., Chang, L.-C., Tsai, M.-J., Wang, Y.-F., & Chang, F.-J. (2010). Dynamic neural
- 1025 networks for real-time water level predictions of sewerage systems-covering gauged and
- 1026 ungauged sites. *Hydrology and Earth System Sciences*, *14*(7), 1309–1319.
- 1027 https://doi.org/10.5194/hess-14-1309-2010
- 1028 DARPA. (2016). Perspective on AI. https://www.darpa.mil/about-us/darpa-perspective-on-ai
- 1029 Dempsey, P., Eadon, A., & Morris, G. (1997). Simpol: A simplified urban pollution modelling
- 1030 tool. Water Science and Technology, 36(8–9), 83–88. https://doi.org/10.1016/S0273-
- 1031 1223(97)00615-X

- 1032 Deuerlein, J. W. (2008). Decomposition Model of a General Water Supply Network Graph.
- 1033 *134*(6), 822–832. https://doi.org/10.1061/(ASCE)0733-9429(2008)134
- 1034 Dini, M., & Tabesh, M. (2017). Water distribution network quality model calibration: A case
- 1035 study-Ahar. *Water Science and Technology: Water Supply*, *17*(3), 759–770.
- 1036 https://doi.org/10.2166/ws.2016.166
- 1037 Dini, M., & Tabesh, M. (2019). Optimal renovation planning of water distribution networks
- 1038 considering hydraulic and quality reliability indices. *Urban Water Journal*, *16*(4), 249–258.
- 1039 https://doi.org/10.1080/1573062X.2019.1669185
- 1040 Doorn, N. (2021). Artificial intelligence in the water domain: Opportunities for responsible use.
- 1041 Science of the Total Environment, 755, 142561.
- 1042 https://doi.org/10.1016/j.scitotenv.2020.142561
- 1043 Dosilovic, F. K., Brcic, M., & Hlupic, N. (2018). Explainable artificial intelligence: A survey.
- 1044 2018 41st International Convention on Information and Communication Technology,
- 1045 *Electronics and Microelectronics, MIPRO 2018 Proceedings, 210–215.*
- 1046 https://doi.org/10.23919/MIPRO.2018.8400040
- 1047 Du, S., Li, T., Yang, Y., & Horng, S. (2020). Neurocomputing Multivariate time series
- 1048 forecasting via attention-based encoder decoder framework. *Neurocomputing*, 388, 269–
- 1049 279. https://doi.org/10.1016/j.neucom.2019.12.118
- Fernandez, G., Livermore, L., Park, C., Kim, N. H., & Haftka, R. (2017). *Review of multi-fidelity models. March.*
- 1052 Fu, G., Kapelan, Z., Kasprzyk, J. R., & Reed, P. (2013). Optimal Design of Water Distribution
- 1053 Systems Using Many-Objective Visual Analytics. *Journal of Water Resources Planning*
- 1054 and Management, 139(6), 624–633. https://doi.org/10.1061/(asce)wr.1943-5452.0000311

- 1055 Gawlikowski, J., Tassi, C. R. N., Ali, M., Lee, J., Humt, M., Feng, J., Kruspe, A., Triebel, R.,
- Jung, P., Roscher, R., Shahzad, M., Yang, W., Bamler, R., & Zhu, X. X. (2021). A Survey of
- 1057 Uncertainty in Deep Neural Networks. 1–41. http://arxiv.org/abs/2107.03342
- 1058 Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press.
- 1059 Hadjimichael, A., Comas, J., & Corominas, L. (2016). Do machine learning methods used in
- 1060 data mining enhance the potential of decision support systems? A review for the urban
- 1061 water sector. AI Communications, 29(6), 747–756. https://doi.org/10.3233/AIC-160714
- 1062 Hang, I., Li, T., Fong, S., & Wong, R. K. (2020). Knowledge-Based Systems Predicting
- 1063 concentration levels of air pollutants by transfer learning and recurrent neural network \Rightarrow .
- 1064 Knowledge-Based Systems, 192, 105622. https://doi.org/10.1016/j.knosys.2020.105622
- Herrera, M., Abraham, E., & Stoianov, I. (2016). A Graph-Theoretic Framework for Assessing
 the Resilience of Sectorised Water Distribution Networks. *Water Resources Management*,
- 1067 *30*(5), 1685–1699. https://doi.org/10.1007/s11269-016-1245-6
- Hoedt, P.-J., Kratzert, F., Klotz, D., Halmich, C., Holzleitner, M., Nearing, G., Hochreiter, S., &
 Klambauer, G. (2021). *MC-LSTM: Mass-Conserving LSTM*.
- 1070 http://arxiv.org/abs/2101.05186
- Huang, C.-L., Hsu, N.-S., Wei, C.-C., & Luo, W.-J. (2015). Optimal spatial design of capacity
 and quantity of rainwater harvesting systems for urban flood mitigation. *Water*
- 1073 (*Switzerland*), 7(9), 5173–5202. https://doi.org/10.3390/w7095173
- 1074 IPCC. (2021). IPCC: Climate Change 2021: The Physical Science Basis. In *Cambridge*
- 1075 University Press. In Press. https://www.ipcc.ch/report/ar6/wg1/
- 1076 IWA. (2021). Digital Water Operational digital twins in the urban water sector: case studies.
- 1077 Jia, Y., Zheng, F., Maier, H. R., Ostfeld, A., Creaco, E., Savic, D., Langeveld, J., & Kapelan, Z.

- 1078 (2021). Water quality modeling in sewer networks: Review and future research directions.
- 1079 *Water Research*, 202(November 2020), 117419.
- 1080 https://doi.org/10.1016/j.watres.2021.117419
- 1081 Kapelan, Z. S., Savic, D. A., & Walters, G. A. (2005). Multiobjective design of water
- 1082 distribution systems under uncertainty. *Water Resources Research*, 41(11), 1–15.
- 1083 https://doi.org/10.1029/2004WR003787
- 1084 Kasprzyk, J. R., Nataraj, S., Reed, P. M., & Lempert, R. J. (2013). Many objective robust
- 1085 decision making for complex environmental systems undergoing change. *Environmental*
- 1086 *Modelling and Software*, 42, 55–71. https://doi.org/10.1016/j.envsoft.2012.12.007
- 1087 Keogh, E., & Mueen, A. (2017). Curse of Dimensionality. *Encyclopedia of Machine Learning*1088 *and Data Mining*, 314–315. https://doi.org/10.1007/978-1-4899-7687-1 192
- 1089 Keum, H. J., Han, K. Y., & Kim, H. I. (2020). Real-Time Flood Disaster Prediction System by
- 1090 Applying Machine Learning Technique. *KSCE Journal of Civil Engineering*, 24(9), 2835–
- 1091 2848. https://doi.org/10.1007/s12205-020-1677-7
- 1092 Kim, H. I., & Han, K. Y. (2020). Urban flood prediction using deep neural network with data
- 1093 augmentation. *Water (Switzerland)*, *12*(3). https://doi.org/10.3390/w12030899
- 1094 Kim, H. I., Keum, H. J., & Han, K. Y. (2019). Real-time urban inundation prediction combining
- 1095 hydraulic and probabilistic methods. *Water (Switzerland)*, *11*(2).
- 1096 https://doi.org/10.3390/w11020293
- 1097 Kratzert, F., Klotz, D., Shalev, G., Klambauer, G., Hochreiter, S., & Nearing, G. (2019).
- 1098 Towards learning universal, regional, and local hydrological behaviors via machine
- *learning applied to large-sample datasets.* 5089–5110.
- 1100 Kumar, A., Rizvi, S. A. A., Brooks, B., Vanderveld, A., Wilson, K. H., Kenney, C., Edelstein, S.,

- Finch, A., Maxwell, A., Zuckerbraun, J., & Ghani, R. (2018). Using Machine Learning to
 Assess the Risk of and Prevent Water Main Breaks. 2657, 1–9.
- 1103 https://doi.org/10.1145/nnnnnnnnn
- 1104 Langeveld, J. G., Benedetti, L., Klein, J. J. M. De, Nopens, I., Amerlinck, Y., Nieuwenhuijzen,
- 1105 V., Flameling, T., Zanten, O. Van, Weijers, S., Benedetti, L., Klein, J. J. M. De, Nopens, I.,
- 1106 & Amerlinck, Y. (2013). Impact-based integrated real-time control for improvement of the
- 1107 Dommel River water quality. 9006. https://doi.org/10.1080/1573062X.2013.820332
- 1108 Latifi, M., Rakhshandehroo, G., Nikoo, M. R., & Sadegh, M. (2019). A game theoretical low
- 1109 impact development optimization model for urban storm water management. *Journal of*
- 1110 *Cleaner Production*, 241. https://doi.org/10.1016/j.jclepro.2019.118323
- Liao, W., Bak-Jensen, B., Pillai, J. R., Wang, Y., & Wang, Y. (2021). A Review of Graph Neural
 Networks and Their Applications in Power Systems. 1–16. http://arxiv.org/abs/2101.10025
- 1113 Lima, G. M., Brentan, B. M., Manzi, D., & Luvizotto, E. (2018). Metamodel for nodal pressure
- 1114 estimation at near real-time in water distribution systems using artificial neural networks.

```
1115 Journal of Hydroinformatics, 20(2), 486–496. https://doi.org/10.2166/hydro.2017.036
```

- 1116 Liu, X., Tian, S., Tao, F., & Yu, W. (2021). Review article A review of artificial neural networks
- in the constitutive modeling of composite materials. *Composites Part B*, 224(May), 109152.
- 1118 https://doi.org/10.1016/j.compositesb.2021.109152
- 1119 Loquercio, A., Segu, M., & Scaramuzza, D. (2020). A General Framework for Uncertainty
- 1120 Estimation in Deep Learning. *IEEE Robotics and Automation Letters*, 5(2), 3153–3160.
- 1121 https://doi.org/10.1109/LRA.2020.2974682
- 1122 Lund, N. S. V., Falk, A. K. V., Borup, M., Madsen, H., & Mikkelsen, P. S. (2018). Model
- 1123 predictive control of urban drainage systems : A review and perspective towards smart real-

- time water management. Critical Reviews in Environmental Science and Technology, 48(3),
- 1125 279–339. https://doi.org/10.1080/10643389.2018.1455484
- 1126 Maier, H., & Dandy, G. (2000). Neural networks for the prediction and forecasting of water
- 1127 resources variables: A review of modelling issues and applications. *Environmental*
- 1128 *Modelling and Software*, 15(1), 101–124. https://doi.org/10.1016/S1364-8152(99)00007-9
- 1129 Maier, H., Kapelan, Z., Kasprzyk, J., Kollat, J., Matott, L. S., Cunha, M. C., Dandy, G. C.,
- 1130 Gibbs, M. S., Keedwell, E., Marchi, A., Ostfeld, A., Savic, D., Solomatine, D. P., Vrugt, J.
- 1131 A., Zecchin, A. C., Minsker, B. S., Barbour, E. J., Kuczera, G., Pasha, F., ... Reed, P. M.
- 1132 (2014). Evolutionary algorithms and other metaheuristics in water resources: Current status,
- research challenges and future directions. *Environmental Modelling and Software*, 62, 271–
- 1134 299. https://doi.org/10.1016/j.envsoft.2014.09.013
- 1135 Martínez, F., Alonso, M., Herna, V., Rao, Z., & Alvisi, S. (2007). Optimizing the operation of
- 1136 the Valencia water- distribution network. 65–78. https://doi.org/10.2166/hydro.2006.018
- 1137 Meijer, D., Bijnen, M. van, Langeveld, J., Korving, H., Post, J., & Clemens, F. (2018).
- 1138 Identifying critical elements in sewer networks using graph-theory. *Water (Switzerland)*,
- 1139 *10*(2). https://doi.org/10.3390/w10020136
- 1140 Meijer, D., Post, J., van der Hoek, J. P., Korving, H., Langeveld, J., & Clemens, F. (2020).
- 1141 Identifying critical elements in drinking water distribution networks using graph theory.
- 1142 *Structure and Infrastructure Engineering*, *17*(3), 347–360.
- 1143 https://doi.org/10.1080/15732479.2020.1751664
- 1144 Meirelles, G., Manzi, D., Brentan, B., Goulart, T., & Luvizotto, E. (2017). Calibration Model for
- 1145 Water Distribution Network Using Pressures Estimated by Artificial Neural Networks.
- 1146 Water Resources Management, 31(13), 4339–4351. https://doi.org/10.1007/s11269-017-

- 1147 1750-2
- Modesto De Souza, R. G., Melo Brentan, B., & Meirelles Lima, G. (2021). *Optimal architecture for artificial neural networks as pressure estimator*. 1–9.
- 1150 Ng, J. Y., Asce, S. M., Fazlollahi, S., Ph, D., Galelli, S., & Asce, M. (2020). Do Design Storms
- 1151 Yield Robust Drainage Systems ? How Rainfall Duration, Intensity, and Profile Can Affect
- 1152 Drainage Performance. 146(3), 1–13. https://doi.org/10.1061/(ASCE)WR.1943-
- 1153 5452.0001167
- 1154 Paluszczyszyn, D., Skworcow, P., & Ulanicki, B. (2013). Online simplification of water
- distribution network models for optimal scheduling. *Journal of Hydroinformatics*, 15(3),
- 1156 652–665. https://doi.org/10.2166/HYDRO.2013.029
- Pascanu, R., Mikolov, T., & Bengio, Y. (2013). On the difficulty of training recurrent neural
 networks. *30th International Conference on Machine Learning, ICML 2013, PART 3*,
- 1159 2347–2355.
- 1160 Pasha, M. F. K., & Lansey, K. (2014). Strategies to develop warm solutions for real-time pump
- scheduling for water distribution systems. *Water Resources Management*, 28(12), 3975–
- 1162 3987. https://doi.org/10.1007/s11269-014-0721-0
- 1163 Paszke, A., Lerer, A., Killeen, T., Antiga, L., Yang, E., Gross, S., Bradbury, J., Massa, F., &
- 1164 Steiner, B. (2019). *PyTorch : An Imperative Style , High-Performance Deep Learning*
- 1165 *Library. NeurIPS.*
- 1166 Pedersen, A. N., Pedersen, J. W., Vigueras-Rodriguez, A., Brink-Kjær, A., Borup, M., &
- 1167 Mikkelsen, P. S. (2021). The Bellinge data set: Open data and models for community-wide
- 1168 urban drainage systems research. *Earth Syst. Sci. Data Discuss.*, *April*, 1–28.
- 1169 https://doi.org/10.5194/essd-2021-8

- 1170 Peng, G. C. Y., Alber, M., Buganza, A., William, T., Suvranu, R. C., Dura, D. S., Garikipati, K.,
- 1171 Karniadakis, G., Lytton, W. W., Perdikaris, P., Petzold, L., & Kuhl, E. (2020). Multiscale
- 1172 Modeling Meets Machine Learning : What Can We Learn? Archives of Computational
- 1173 *Methods in Engineering*, 0123456789. https://doi.org/10.1007/s11831-020-09405-5
- 1174 Prasad, T. D., & Park, N.-S. (2004). Multiobjective Genetic Algorithms for Design of Water
- 1175 Distribution Networks. Journal of Water Resources Planning and Management, 130(1), 73–
- 1176 82. https://doi.org/10.1061/(asce)0733-9496(2004)130:1(73)
- 1177 Raei, E., Reza Alizadeh, M., Reza Nikoo, M., & Adamowski, J. (2019). Multi-objective
- decision-making for green infrastructure planning (LID-BMPs) in urban storm water
- 1179 management under uncertainty. *Journal of Hydrology*, 579.
- 1180 https://doi.org/10.1016/j.jhydrol.2019.124091
- 1181 Rao, Z., & Alvarruiz, F. (2007). Use of an artificial neural network to capture the domain
- 1182 *knowledge of a conventional hydraulic simulation model*. 15–24.
- 1183 https://doi.org/10.2166/hydro.2006.014
- 1184 Rao, Z., & Salomons, E. (2007). Development of a real-time, near-optimal control process for
- 1185 *water-distribution networks Zhengfu Rao and Elad Salomons.* 25–37.
- 1186 https://doi.org/10.2166/hydro.2006.015
- 1187 Razavi, S., Jakeman, A., Saltelli, A., Prieur, C., Iooss, B., Borgonovo, E., Plischke, E., Lo Piano,
- 1188 S., Iwanaga, T., Becker, W., Tarantola, S., Guillaume, J. H. A., Jakeman, J., Gupta, H.,
- 1189 Melillo, N., Rabitti, G., Chabridon, V., Duan, Q., Sun, X., ... Maier, H. R. (2021). The
- 1190 Future of Sensitivity Analysis: An essential discipline for systems modeling and policy
- support. *Environmental Modelling and Software*, *137*(December 2020).
- 1192 https://doi.org/10.1016/j.envsoft.2020.104954

- 1193 Razavi, S., Tolson, B. A., & Burn, D. H. (2012a). Numerical assessment of metamodelling
- 1194 strategies in computationally intensive optimization. *Environmental Modelling and*

1195 Software, 34, 67–86. https://doi.org/10.1016/j.envsoft.2011.09.010

- 1196 Razavi, S., Tolson, B. A., & Burn, D. H. (2012b). Review of surrogate modeling in water
- resources. Water Resources Research, 48(7). https://doi.org/10.1029/2011WR011527
- 1198 Roscher, R., Bohn, B., Duarte, M., & Garcke, J. (2020). Explainable Machine Learning for

1199 Scientific Insights and Discoveries. *IEEE Access*, 8, 42200–42216.

- 1200 https://doi.org/10.1109/ACCESS.2020.2976199
- 1201 Rosin, T. R., Romano, M., Keedwell, E., & Kapelan, Z. (2021). A Committee Evolutionary
- 1202 Neural Network for the Prediction of Combined Sewer Overflows. *Water Resources*

1203 *Management*, 35(4), 1273–1289. https://doi.org/10.1007/s11269-021-02780-z

1204 Salomons, E., Goryashko, A., Shamir, U., Rao, Z., & Alvisi, S. (2007). Optimizing the operation

1205 of the Haifa-A water-distribution network. 51–64. https://doi.org/10.2166/hydro.2006.017

- 1206 Sambito, M., Di Cristo, C., Freni, G., & Leopardi, A. (2020). Optimal water quality sensor
- 1207 positioning in urban drainage systems for illicit intrusion identification. *Journal of*

1208 *Hydroinformatics*, 22(1), 46–60. https://doi.org/10.2166/hydro.2019.036

- 1209 Sayers, W., Savic, D., & Kapelan, Z. (2019). Performance of LEMMO with artificial neural
- networks for water systems optimisation. *Urban Water Journal*, *16*(1), 21–32.
- 1211 https://doi.org/10.1080/1573062X.2019.1611886
- 1212 Schultz, M. T., Small, M. J., Farrow, R. S., & Fischbeck, P. S. (2004). State Water Pollution
- 1213 Control Policy Insights from a Reduced-Form Model. *Journal of Water Resources Planning*
- 1214 and Management, 130(2), 150–159. https://doi.org/10.1061/(ASCE)0733-
- 1215 9496(2004)130:2(150)

- 1216 Seyedashraf, O., Bottacin-Busolin, A., & Harou, J. J. (2021). A Disaggregation-Emulation
- 1217 Approach for Optimization of Large Urban Drainage Systems Water Resources Research.
- 1218 2017, 1–18. https://doi.org/10.1029/2020WR029098
- 1219 Shamir, U., Asce, F., & Salomons, E. (2008). Optimal Real-Time Operation of Urban Water
- 1220 Distribution Systems Using Reduced Models. *Journal of Water Resources Planning and*
- 1221 *Management*, 134(2), 181–185. https://doi.org/10.1061/(ASCE)0733-9496(2008)134:2(181)
- 1222 She, L., & You, X.-Y. (2019). A Dynamic Flow Forecast Model for Urban Drainage Using the
- 1223 Coupled Artificial Neural Network. *Water Resources Management*, *33*(9), 3143–3153.
- 1224 https://doi.org/10.1007/s11269-019-02294-9
- 1225 Shen, C. (2018). A Transdisciplinary Review of Deep Learning Research and Its Relevance for
- 1226 Water Resources Scientists. *Water Resources Research*, 54(11), 8558–8593.
- 1227 https://doi.org/10.1029/2018WR022643
- 1228 Snyder, H. (2019). Literature review as a research methodology: An overview and guidelines.
- 1229 *Journal of Business Research*, 104(March), 333–339.
- 1230 https://doi.org/10.1016/j.jbusres.2019.07.039
- 1231 Strubell, E., Ganesh, A., & McCallum, A. (2020). Energy and policy considerations for modern
- deep learning research. AAAI 2020 34th AAAI Conference on Artificial Intelligence, 1,
- 1233 1393–13696. https://doi.org/10.1609/aaai.v34i09.7123
- 1234 Sun, Y., Wang, Y., Fu, K., Wang, Z., Zhang, C., & Ye, J. (2020). Constructing geographic and
- 1235 long-term temporal graph for traffic forecasting. *Proceedings International Conference on*
- 1236 Pattern Recognition, 3483–3490. https://doi.org/10.1109/ICPR48806.2021.9412506
- 1237 Tscheikner-Gratl, F., Zeisl, P., Kinzel, C., Rauch, W., Kleidorfer, M., Leimgruber, J., & Ertl, T.
- 1238 (2016). Lost in calibration: Why people still do not calibrate their models, and why they still

- should A case study from urban drainage modelling. *Water Science and Technology*,
- 1240 74(10), 2337–2348. https://doi.org/10.2166/wst.2016.395
- 1241 Tsiami, L., & Makropoulos, C. (2021). Cyber—physical attack detection in water distribution
- systems with temporal graph convolutional neural networks. *Water (Switzerland)*, *13*(9).
- 1243 https://doi.org/10.3390/w13091247
- Ulanicki, B., Zehnpfund, A., & Martinez, F. (1996). Simplification of Water Distribution
 Network Models. September. https://doi.org/10.13140/RG.2.1.4340.8404
- 1246 UNESCO. (2020). United Nations World Water Development Report 2020: Water and Climate
 1247 Change.
- University of Kentucky. (2013). Water Distribution System Research Database, University of
 Kentucky. https://doi.org/10.13023/kwrri.wdsrd.
- 1250 van Bijnen, M., Korving, H., Langeveld, J., & Clemens, F. (2017). Calibration of hydrodynamic
- 1251 model-driven sewer maintenance. *Structure and Infrastructure Engineering*, 13(9), 1167–
- 1252 1185. https://doi.org/10.1080/15732479.2016.1247287
- Verma, S., & Zhang, Z.-L. (2019). Learning Universal Graph Neural Network Embeddings With
 Aid Of Transfer Learning.
- 1255 Vogado, L. H. S., Veras, R. M. S., Araujo, F. H. D., Silva, R. R. V, & Aires, R. T. (2018).
- 1256 Engineering Applications of Artificial Intelligence Leukemia diagnosis in blood slides using
- 1257 transfer learning in CNNs and SVM for classification. *Engineering Applications of*
- 1258 Artificial Intelligence, 72(October 2017), 415–422.
- 1259 https://doi.org/10.1016/j.engappai.2018.04.024
- 1260 Vojinovic, Z., Kecman, V., & Babovic, V. (2003). Hybrid approach for modeling wet weather
- response in wastewater systems. Journal of Water Resources Planning and Management,

1262	129(6), 511–521. ht	tps://doi.org/10.106	51/(ASCE)0733-94	96(2003)129:6(511)
			· · · · · ·	

- 1263 Wang, Q., Guidolin, M., Savic, D., & Kapelan, Z. (2015). Two-Objective Design of Benchmark
- 1264 Problems of a Water Distribution System via MOEAs: Towards the Best-Known
- 1265 Approximation of the True Pareto Front. Journal of Water Resources Planning and
- 1266 *Management*, 141(3), 04014060. https://doi.org/10.1061/(asce)wr.1943-5452.0000460
- 1267 Weiss, K., Khoshgoftaar, T. M., & Wang, D. (2016). A survey of transfer learning. In *Journal of*
- 1268 *Big Data*. Springer International Publishing. https://doi.org/10.1186/s40537-016-0043-6
- 1269 Wolfs, V., & Willems, P. (2017). Modular Conceptual Modelling Approach and Software for
- 1270 Sewer Hydraulic Computations. *Water Resources Management*, *31*(1), 283–298.
- 1271 https://doi.org/10.1007/s11269-016-1524-2
- Wong, E., & Kolter, J. Z. (2019). *Learning perturbation sets for robust machine learning*. 1–32.
 https://arxiv.org/pdf/2007.08450.pdf
- 1274 Wu, L., Zulueta, K., Major, Z., Arriaga, A., & Noels, L. (2020). Bayesian inference of non-linear
- 1275 multiscale model parameters accelerated by a Deep Neural Network. *Computer Methods in*
- 1276 *Applied Mechanics and Engineering*, 360, 112693.
- 1277 https://doi.org/10.1016/j.cma.2019.112693
- 1278 Ying, R., Bourgeois, D., You, J., Zitnik, M., & Leskovec, J. (2019). GNNExplainer: Generating

1279 explanations for graph neural networks. *Advances in Neural Information Processing*1280 *Systems*, *32*(iii).

- 1281 Yoon, S., Lee, Y.-J., & Jung, H.-J. (2020). Accelerated monte carlo analysis of flow-based
- 1282 system reliability through artificial neural network-based surrogate models. *Smart*
- 1283 Structures and Systems, 26(2), 175–184. https://doi.org/10.12989/sss.2020.26.2.175
- 1284 Zhang, D., Mishra, S., Brynjolfsson, E., Etchemendy, J., Ganguli, D., Grosz, B., Lyons, T.,

- 1285 Manyika, J., Niebles, J. C., Sellitto, M., Shoham, Y., Clark, J., & Perrault, R. (2021). 2021
- 1286 AI Index Report. 1–222. https://aiindex.stanford.edu/report/
- 1287 Zhang, W., Li, J., Chen, Y., & Li, Y. (2019). A Surrogate-Based Optimization Design and
- 1288 Uncertainty Analysis for Urban Flood Mitigation. *Water Resources Management*, 33(12),
- 1289 4201–4214. https://doi.org/10.1007/s11269-019-02355-z
- 1290 Zhou, J., Cui, G., Zhang, Z., Yang, C., Liu, Z., Wang, L., Li, C., & Sun, M. (2018). Graph
- 1291 *Neural Networks: A Review of Methods and Applications.* 1–22.
- 1292 Zhu, Y., & Zabaras, N. (2018). Bayesian deep convolutional encoder decoder networks for
- surrogate modeling and uncertainty quantification. *Journal of Computational Physics*, 366,
- 1294 415–447. https://doi.org/10.1016/j.jcp.2018.04.018