

This document is the Accepted Manuscript version of a Published Work that appeared in final form in ACS ES&T Water, copyright © 2025 American Chemical Society after peer review and technical editing by the publisher. To access the final edited and published work see <https://doi.org/10.1021/acsestwater.4c01220>.

# 1     **Modeling of a Mainstream Partial Nitrification/Anammox Process through a Hybrid** 2                                    **Theoretical-Machine Learning Approach**

3     Valeria Alvarado<sup>1†</sup>, Lebing Ying<sup>1†</sup>, Vahid Asghari<sup>1</sup>, Shu-Chien Hsu<sup>1\*</sup>, and Po-Heng Lee<sup>2</sup>

4                   <sup>1</sup> Department of Civil and Environmental Engineering, The Hong Kong Polytechnic University, Hong Kong S.A.R

5                   <sup>2</sup> Department of Civil and Environmental Engineering, Imperial College London, South Kensington SW7 2AZ, London, United Kingdom

6     †These authors contributed equally to this work and were designated as co-first authors.

7     \* Corresponding Author: [mark.hsu@polyu.edu.hk](mailto:mark.hsu@polyu.edu.hk), Tel: +852-27666057

8

9     **Abstract:** Model simulations are vital in optimizing and predicting the performance of biological  
10  sewage treatment, especially for processes involving slow-growing bacteria. However, data records  
11  often include missing, invalid, or infrequent measurements of parameters, compromising prediction  
12  accuracy. This study used a hybrid theoretical-machine learning approach to address these issues. By  
13  leveraging stoichiometry and kinetics, missing values were calculated in limited data sets, which were  
14  then analyzed through machine learning algorithms to reveal hidden microbial interactions. The model  
15  was validated with data from a pilot-scale partial nitritation/anammox fluidized bed membrane  
16  bioreactor (PN/A FMBR) with saline sewage. The model demonstrated strong prediction performance,  
17  with random forest outperforming other algorithms, with correlation coefficients of 0.89, 0.72, and 0.80  
18  for ammonium, nitrite, and nitrate data sets, respectively, when compared to actual values. Training sets  
19  containing 73 to 88 same-day values reached acceptable predicting performance. The results also  
20  revealed that microbial synergy in nitrogen transformation, particularly in the partial denitrification  
21  from nitrate to nitrite linked to Anammox in responding to a low DO supply, was evident in this PN/A  
22  FMBR. Additionally, key parameters, including temperature, pH, and specific microbiomes, were  
23  identified as critical for predicting PN/AFMBR performance, highlighting significant microbial  
24  interactions that warrant further investigation.

25     **Keywords:** Wastewater treatment, Anammox, Theoretical-machine learning, Partial nitritation,  
26     Microbial interactions, Fluidized bed membrane bioreactor

27     **Synopsis:** This study presents an innovative hybrid Theoretical-Machine Learning model for sewage  
28     treatment, calculating missing values and revealing crucial microbial interactions to enhance effluent  
29     prediction accuracy.

30

## 31 1. INTRODUCTION

32 Traditional biological nitrogen removal (BNR) methods in sewage treatment plants (STPs),  
33 such as nitrification/denitrification (N/D) are energy-intensive and produce high levels of  
34 sludge and green house gas. Alternative approaches like partial nitrification (PN) and anaerobic  
35 ammonium oxidation (anammox) offer more sustainable solutions.<sup>1</sup> PN nitrifies half the  
36 ammonium ( $\text{NH}_4^+$ ) in the influent to nitrite ( $\text{NO}_2^-$ ) by ammonium oxidizing bacteria, while  
37 anammox removes the remaining  $\text{NH}_4^+$  with  $\text{NO}_2^-$  to produce nitrogen gas ( $\text{N}_2$ ) by anammox  
38 bacteria<sup>2-7</sup>. Coupled PN/anammox (PN/A) process produces  $\text{N}_2$  and small amounts of nitrate  
39 ( $\text{NO}_3^-$ ). Integration of the PN/A process in a membrane bioreactor (MBR) has shown high  
40 nitrogen removal efficiency (~89.5%) has been theoretically demonstrated<sup>3,8</sup>. In addition to  
41 the benefits of the PN/A processes, its integration to an MBR results in high-quality effluent<sup>9-</sup>  
42 <sup>11</sup> at short hydraulic retention times (HRTs) and high loading rates.<sup>2,12</sup>

43 Research indicates that the PN/A process is complicated and highly susceptible to  
44 environmental factors, including temperature, pH, dissolved oxygen, the presence of heavy  
45 metals and the concentration of ammonium and nitrite<sup>13</sup>. Additionally, the sensitivity of  
46 anammox bacteria and the unstable suppression of nitrite-oxidizing bacteria (NOB) in the  
47 complex mainstream environment vastly limited anammox applications.<sup>14,15</sup> This sensitivity  
48 necessitates precise control of operational parameters to avoid costly and time-consuming  
49 recovery from disturbances..<sup>16-18</sup> Modelling has proven to be an invaluable tool for enhancing  
50 understanding of recent anammox-based processes and helping in the simulating various  
51 operational strategies<sup>19-21</sup>. In the context of PN/A systems, the primary focus was on  
52 autotrophic nitrogen removal carried out by ammonia-oxidizing bacteria (AOB), NOB and  
53 anammox bacteria. However, recent research suggests that heterotrophic microorganisms may  
54 also play a critical role in these systems.<sup>20</sup> While autotrophic processes dominate the PN/A  
55 pathway, heterotrophs can contribute by influencing nitrite availability and competing for

56 resources, thereby affecting overall system efficiency. Therefore, more efforts are required to  
57 characterize the microbial interactions not only between AOB and NOB but also between  
58 autotrophic and heterotrophic organisms to improve nitrogen removal efficiency.

59 Understanding these microbial synergies, alongside the multiple pathways contributing to  
60 nitrogen transformation, is crucial for optimizing PN/A performance. Therefore, further  
61 research is necessary to better characterize microbial interactions by identifying multiple  
62 pathway contributions, influencing factors and energy-efficient operational strategies for the  
63 successful implementation of these processes.<sup>19,22</sup>

64 Models which include microbial community diversity enable the understanding of microbial  
65 community shifts and composition link to control strategies, stability, and performance.<sup>23</sup> To  
66 address the complexities of modeling such systems, tools like the activated sludge model  
67 (ASM) series<sup>24,25</sup> and anaerobic digestion model 1,<sup>26</sup> are commonly used, integrated with  
68 software such as AQUASIM<sup>27</sup>, GPS-X<sup>28</sup>, and Biowin<sup>29</sup> with special modules for MBR  
69 simulations<sup>30,31</sup>. PN/A process in granular sludge reactors have been studied with AQUASIM's  
70 multi-substrate<sup>17,32-35</sup> and multi-species biofilm models. These models are invaluable for  
71 predicting system behavior and optimizing operational strategies<sup>36</sup> but often suffer from a lack  
72 of standardized methods for parameter calculation and challenges in scaling up from biofilm-  
73 based models to real-world applications.<sup>32,37</sup>

74 In recent years, the statistical-empirical approach through machine learning algorithms have  
75 been applied in the sewage treatment field for sewage quality influent and effluent prediction<sup>38-</sup>  
76 <sup>40</sup>, soft measurement<sup>41</sup> and automatic control of STPs<sup>42,43</sup> due to its ability for handling large  
77 datasets and dealing with complex, nonlinear relationships and accurate predicting.  
78 Furthermore, recent advances in interpretable techniques,<sup>44,45</sup> now offer insights into the key  
79 factors driving biological sewage treatment. In practical wastewater treatment scenarios, data  
80 gaps are common,<sup>37,46</sup> particularly with key operational parameters like microbial community

81 abundances and real-time nitrogen concentrations.<sup>45,47</sup> These gaps can significantly hinder the  
82 accuracy of outcome predictions. Additionally, relying solely on machine learning can fail to  
83 capture the complex biological interactions within these systems, often resulting in incomplete  
84 process understanding.<sup>48</sup> While there have been many models developed for PN/A process<sup>19,86</sup>,  
85 few theoretical- or statistical-empirical modelling approaches have specifically addressed the  
86 PN/A process in an anoxic MBR with granulated activated carbon (GAC) and a fluidized-bed  
87 with sewage recirculation. Moreover, the influence of microbial community abundance on  
88 nitrogen removal in this distinctive setup has seldom been analyzed through theoretical or  
89 statistical-empirical modeling approach. In this context, this paper introduces a novel hybrid  
90 theoretical-statistical approach designed specifically for data-limited settings. This approach is  
91 applied to a partial nitrification/anammox fluidized-bed membrane bioreactor (PN/AFMBR)  
92 to predict effluent nitrogen compounds effectively. It also aims to identify crucial parameters  
93 influencing biological nitrogen removal (BNR) by incorporating operational, water quality, and  
94 microbial data. By calculating missing values and revealing hidden microbial interactions that  
95 traditional methods overlook, this study successfully bridges gaps in optimizing the  
96 performance of biological sewage treatment systems, particularly those relying on slow-  
97 growing bacteria.

## 98 **2. MATERIALS AND METHODS**

### 99 **2.1 PN/AFMBR description**

100 The pilot-scale PN/AFMBR for mainstream treatment of settled sewage obtained from  
101 chemically enhanced primary treatment was installed at the Stonecutters Island sewage  
102 treatment works (STWs) in Hong Kong, as shown in Figure S1<sup>49</sup>. Sewage in Hong Kong  
103 contains high salinity because about 80% of the population is served with treated seawater for  
104 toilet flushing. The reactor operated for 523 days with final working volume of 1.23 m<sup>3</sup><sup>4</sup> The

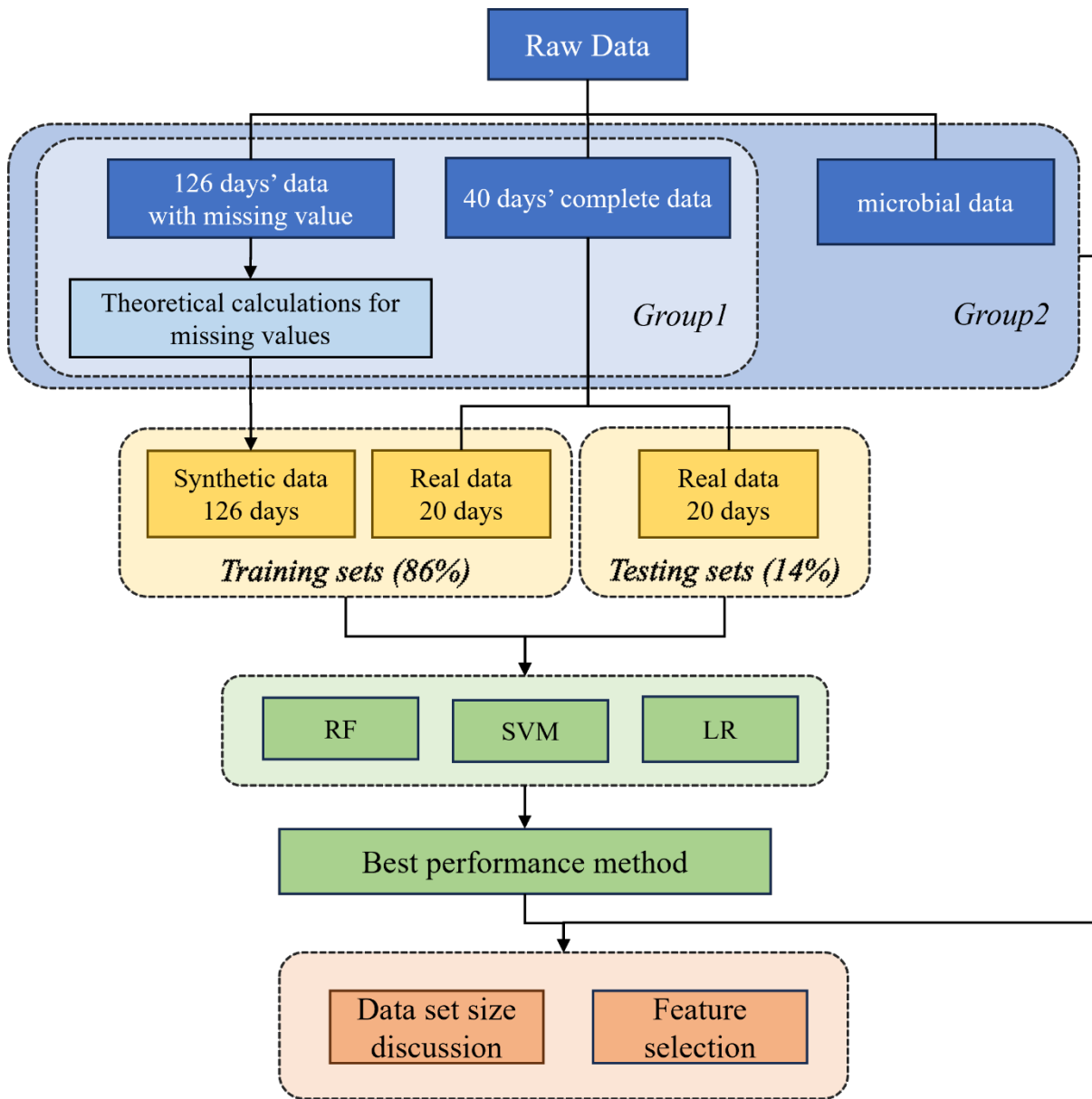
105 acclimation period was about 147 days when the operation was shifted from batch to  
106 continuous mode. Settled GAC occupied 56% of the reactor, while fluidized GAC occupied  
107 100% of the reactor height. The PN/AFMBR contained five submerged hollow-fiber  
108 membrane (pore size: 0.1  $\mu\text{m}$ ) modules made of polyvinylidene fluoride with a total surface  
109 area of 55  $\text{m}^2$ . The upflow velocity was 69.2 m/h. The reactor was inoculated with seed from  
110 an anaerobic digester in Tai Po STWs.

111 The operation of the PNAFMBR was automated with an integrated programmable logic  
112 controller (PLC) unit including sensors for temperature, dissolved oxygen (DO), and pH.  
113 Transmembrane pressure (TMP) was recorded with a pressure gauge installed in the effluent  
114 line. The average flow rate was adjusted according to three different HRTs (i.e., 15,10 and 8  
115 hr).<sup>4</sup> Chemical oxygen demand (COD), ammonium, nitrite, and nitrate concentrations were  
116 determined according to standard methods.<sup>50</sup>

117 DNA extraction of biomass samples (i.e., seed sludge, suspension sludge in the PN/AFMBR,  
118 and biofilm from GAC in the PN/AFMBR) was performed with PowerSoil DNA extraction kit  
119 (MoBio Laboratories, Carlsbad, CA, USA). 16S rRNA gene was sequenced with Illumina Mi-  
120 Seq platform (PE250) using primers 515F and 806R. BlastKOALA was utilized for genome  
121 functions characterization and functional pathways reconstruction<sup>4,51</sup>. Details of the pilot-scale  
122 PN/AFMBR, its operating conditions and microbial community results are reported in Huang  
123 et al<sup>4</sup>.

124 The methodology, as described in **Figure 1**, starts with analyzing data for same-day  
125 information. Missing values are then calculated theoretically. Data sets are separated in two  
126 groups depending on the number of input process parameters. Group 1 includes operation  
127 parameters and influent nitrogen compounds concentrations. Besides the input process  
128 parameters in group 1, microbial community abundances are included in Group 2. The machine  
129 learning algorithms are then developed for each data set in Group 1. The best machine learning

130 algorithm is identified through a ranking approach. Analysis of different data set sizes –  
 131 including Group 2 – and feature selection are conducted with the machine learning algorithm  
 132 with the best performance.



133  
 134  
 135  
 136

**Figure 1.** Overall framework of the hybrid theoretical-statistical approach, RF: random forest, and SVM: support vector machine, LR: linear regression.

## 137 2.2 Dataset preparation

### 138 2.2.1 Operational data description

139 To build the datasets, data was screened for same-day information on six operational  
140 parameters (i.e., HRT, temperature, transmembrane pressure, average flow rate, DO, pH), three  
141 influent and three effluent nitrogen compounds concentrations (i.e., ammonium, nitrite, and  
142 nitrate), and 56 microbial community abundance percentages in the suspension and biofilm  
143 samples in the PN/AFMBR (e.g., *Nitrosomonas*, *Nitrospira*, and *Candidatus Kuenenia*). The  
144 available data for the analysis was collected for 166 days during the period of April 2018 to  
145 October 2018 (day 181 to 346).

146 Excluding temperature, all operational parameters were controlled with the PLC unit.  
147 Temperature, TMP and pH in the reactor ranged from 25 to 35°C, -0.1 to -0.01 kg/cm<sup>2</sup> and 7.2  
148 to 7.8, respectively. pH was increased with sodium carbonate solution when the value  
149 decreased below 7.2. Even if DO set point was 0.1 mg/L, the average DO value from the  
150 measurements of five DO sensors installed along the reactor were recorded. Operational  
151 parameters measurements were frequently recorded (3 to 6 times per week) from the PLC unit.  
152 Concentrations of inorganic nitrogen compounds were routinely measured every 3 to 5 days at  
153 the influent and effluent of the PN/AFMBR. The nitrogen compounds influent concentrations  
154 were considered as uncontrollable variables. Samples for microbial community analysis were  
155 collected 5 times in the total period of data collection on Days 0, 216, 259, 308, and 332.<sup>4</sup> The  
156 explanation of dates and frequency of microbial community analysis sampling is explained in  
157 **Supporting Method S2**. At genus level, 28 microbiomes were identified for each of the seed  
158 sludge sample (S0), suspension sludge samples (S1 to S4), and biofilm samples from the GAC  
159 (B1 to B4). Microbial community abundance percentages are assumed as controllable because  
160 operation parameters might inhibit the growth of a group of microbiomes. For example, limited  
161 DO (< 0.1 mg/L) is expected to suppress nitrite oxidizing bacteria. Yet, the microbial  
162 abundance of different microbiomes in the seed sludge was out of the control of this study.

163 Given the disparity on the frequency of data collection, only one day (Day 332) has same-day  
 164 information for operational parameters, influent and effluent nitrogen concentrations, and  
 165 microbial community abundance. From Day 181 to 346, data for operational parameters were  
 166 collected for 105 days, whereas data for influent and effluent nitrogen compounds were  
 167 measured for 50 days. Data for microbial community abundance was obtained for 5 days.  
 168 Records with missing values have been eliminated<sup>41,52</sup> or calculated from previous data<sup>53</sup> in  
 169 other models based on machine learning methods. In this study, missing data for input and  
 170 output process variables were calculated under theoretical assumptions to provide more  
 171 professional insights.

172 **Table 1** Statistical description of input and output process variables in data sets of Groups 1 and 2 according to  
 173 Huang et al. <sup>4</sup>

Group 1	Group 2	Process variables	Input	Min.	Max.	Mean	Std. Dev.
✓	✓	Operational parameters	HRT (h)	8.000	15.000	12.663	2.881
			Temperature (°C)	24.810	34.46	30.377	1.890
			TMP (kg/cm <sup>2</sup> )	-0.100	-0.010	-0.048	0.016
			Flow rate (m <sup>3</sup> /h)	0.020	0.940	0.109	0.093
			Dissolved oxygen (mg/L)	0.030	0.330	0.090	0.042
✓	✓	Influent (mg/L)	pH	7.160	7.770	7.528	0.070
			Ammonium	11.170	44.000	22.040	5.083
			Nitrite	0.000	1.270	0.025	0.120
	✓	Microbial data	Nitrate	0.000	7.130	2.831	1.629
			Microbial community abundance (%):	*	*	*	*
1	2		Output				
✓	✓	Effluent & influent difference(mg/L)	Ammonium	5.272	43.910	18.809	5.589
			Nitrite	-9.779	1.230	-1.713	2.222
			Nitrate	-14.010	5.470	-1.440	4.278

Notes: \*The statistical description of the microbial community abundance includes 28 genera for suspensions in the PN/AFMBR and 28 genera for biofilm in the GAC of the PN/AFMBR are listed in **Table S8-S9**.

174

175 Input process variables include operational parameters, nitrogen compounds concentrations in  
 176 the influent and microbial community abundance percentages in the suspension and biofilm  
 177 within the PN/AFMBR (**Table 1**). Missing data for TMP and DO were assumed as the average  
 178 between the previous and following day but circumstances on-site were also considered.<sup>54</sup> For  
 179 example, TMP was assumed lower than average during days 239 to 241 due to power shortcuts



180 on-site. On Day 256, the DO was set to 0 mg/L thus the average DO in the reactors was assumed  
181 as 0.04 mg/L as measured for Day 255. Interpolation was used to assume the missing data for  
182 flow rate and pH.<sup>55</sup> In addition, the on-site observations were considered. For example, flow  
183 rate on Day 276 was below the measurement range, thus 0.055 m<sup>3</sup>/hr was assumed. The missing  
184 data for temperature were assumed by the average of the previous and the following day. If  
185 temperature data for two days were missing, then the missing values were interpolated from  
186 one day before and one day after. Missing values for influent ammonium concentration were  
187 assumed as a random value within a range obtained from measurements on April 2018 to  
188 October 2018<sup>4</sup> and April 2016 to October 2016.<sup>56</sup> Nitrite concentrations in the influent were  
189 assumed as 0 mg/L for all the missing data because all the measurements in this period were  
190 below 0.6 mg/L, except for day 274 (1.27 mg/L). Missing data for nitrate concentration in the  
191 influent were assumed as a random number within the range of the measured maximum and  
192 minimum concentrations for each month (April to October 2018). Theoretical calculations  
193 based on stoichiometry and microbial kinetics have been previously determined to describe the  
194 biochemical reactions occurring in the pilot-scale PN/AFMBR through a thermodynamics  
195 electron equivalent (TEE) model.<sup>22</sup> Even though microbial growth can be established through  
196 stoichiometry and kinetics, the microbial community abundance are only limited within known  
197 microbiomes. Given that quantification of total active biomass is required to calculate the  
198 percentage of microbial community abundance for each genus of known microbiomes, missing  
199 percentages were assumed with linear or exponential regression between the 5 days measured.

200 Output process variables were calculated as the difference between effluent and influent  
201 nitrogen compounds concentrations (**Table 1**). Effluent nitrogen compounds concentrations  
202 were measured for 50 days; thus, the TEE model developed by Alvarado<sup>22</sup> was used to validate  
203 the data of 50 measurements from Huang et al.<sup>4</sup>

204 2.2.2 Theoretical calculation of missing data

205 Missing values for effluent concentrations were determined based on the TEE model. The TEE  
 206 model assumes that biological reactions are based on substrate partitioning and cellular  
 207 framework, quantifies microbial community abundances in various sludge types, and includes  
 208 multiple biological processes (e.g., organic oxidation, denitrification, sulfate reduction) with  
 209 iterative calculations for influent and effluent nitrogen compounds, utilizing design criteria  
 210 from existing literature and experimental data to estimate removal efficiencies and account for  
 211 missing data using specified thresholds and regression methods. The detailed calculation  
 212 processes and its correlated assumptions are listed in **Supporting Method S1. Table S1** lists  
 213 typical design criteria for the biological processes in the PN/AFMBR included in the TEE  
 214 model. Ammonium concentration (in mg/L) in the effluent ( $NH_4^+_{eff}$ ) is calculated by including  
 215 ammonium concentration in the influent ( $NH_4^+_{inf}$ ); ammonium produced during organic  
 216 oxidation ( $NH_4^+_{O.O,p}$ ), fermentation ( $NH_4^+_{Fer,p}$ ), heterotrophic denitrification with nitrite  
 217 ( $+NH_4^+_{HD,NO_2^-,p}$ ), heterotrophic partial denitrification with nitrate reduction to nitrite  
 218 ( $NH_4^+_{HPD,NO_3^-,p}$ ), dissimilatory nitrate reduction to ammonium with sulfur reduction  
 219 ( $NH_4^+_{DNRA,p}$ ), sulfate reduction ( $NH_4^+_{SR,p}$ ), and heterotrophic denitrification with nitrate  
 220 ( $NH_4^+_{HD,NO_3^-,p}$ ); and, ammonium consumed during aerobic ammonium oxidation ( $NH_4^+_{AAO,c}$ ),  
 221 nitrite oxidation ( $NH_4^+_{N.O,c}$ ), autotrophic partial denitrification with nitrate reduction to nitric  
 222 oxide ( $NH_4^+_{APD,c}$ ), anammox ( $NH_4^+_{Anx,c}$ ), autotrophic denitrification with nitrate ( $NH_4^+_{AD,NO_3^-,c}$ ),  
 223 and methanol oxidation ( $NH_4^+_{MO,c}$ ) as shown in Eq. 1. Nitrite concentration (in mg/L) in the  
 224 effluent ( $NO_2^-_{eff}$ ) is determined by including nitrite concentration in the influent ( $NO_2^-_{inf}$ );  
 225 nitrite produced in aerobic ammonium oxidation ( $NO_2^-_{AAO,p}$ ), and heterotrophic partial  
 226 denitrification with nitrate reduction to nitrite ( $NO_2^-_{HPD,NO_3^-,p}$ ); and, consumed in heterotrophic  
 227 denitrification with nitrite ( $NO_2^-_{HD,NO_2^-,c}$ ), nitrite oxidation ( $NO_2^-_{N.O,c}$ ), and anammox

228 ( $NO_2^-_{Anx,c}$ ) as shown in Eq.2. Nitrate concentration (in mg/L) in the effluent ( $NO_3^-_{eff}$ ) is  
229 calculated by including nitrate concentration in the influent ( $NO_3^-_{inf}$ ); nitrate produced during  
230 nitrite oxidation ( $NO_3^-_{N.O.,p}$ ), and anammox ( $NO_3^-_{Anx,p}$ ); and, nitrate consumed in autotrophic  
231 partial denitrification with nitrate reduction to nitric oxide ( $NO_3^-_{APD,c}$ ), autotrophic  
232 denitrification with nitrate ( $NO_3^-_{AD,NO_3,c}$ ), dissimilatory nitrate reduction to ammonium with  
233 sulfur reduction ( $NO_3^-_{DNRA,c}$ ), heterotrophic partial denitrification with nitrate reduction to  
234 nitrite ( $NO_3^-_{HPD,NO_3,c}$ ), heterotrophic denitrification with nitrate ( $NO_3^-_{HD,NO_3,c}$ ), methanol  
235 oxidation ( $NO_3^-_{MO,c}$ ), and heterotrophic partial denitrification with nitrate reduction to nitric  
236 oxide by *Lewinellaceae* ( $NO_3^-_{Lew,c}$ ) as shown in Eq.3.

$$237 \quad NH_4^+_{eff} = NH_4^+_{inf} + NH_4^+_{O.O,p} + NH_4^+_{Fer.,p} + NH_4^+_{HD,NO_2,p} + NH_4^+_{HPD,NO_3,p} + NH_4^+_{DNRA,p} +$$

$$238 \quad NH_4^+_{SR,p} + NH_4^+_{HD,NO_3,p} - NH_4^+_{AAO,c} - NH_4^+_{N.O,c} - NH_4^+_{APD,c} - NH_4^+_{Anx,c} - NH_4^+_{AD,NO_3,c} -$$

$$239 \quad NH_4^+_{MO,c} \quad (Eq.1)$$

$$240 \quad NO_2^-_{eff} = NO_2^-_{inf} + NO_2^-_{AAO,p} + NO_2^-_{HPD,NO_3,p} - NO_2^-_{HD,NO_2,c} - NO_2^-_{N.O,c} - NO_2^-_{Anx,c} \quad (Eq.2)$$

$$241 \quad NO_3^-_{eff} = NO_3^-_{inf} + NO_3^-_{N.O.,p} + NO_3^-_{Anx,p} - NO_3^-_{APD,c} - NO_3^-_{AD,NO_3,c} - NO_3^-_{DNRA,c} -$$

$$242 \quad NO_3^-_{HPD,NO_3,c} - NO_3^-_{HD,NO_3,c} - NO_3^-_{MO,c} - NO_3^-_{Lew,c} \quad (Eq.3)$$

243  
244 Operational parameters and influent nitrogen compounds concentrations were the input process  
245 variables in data sets of Group 1, while data sets in Group 2 also included microbial community  
246 abundance percentages as input process variables to determine the influence of microbiome on  
247 the prediction of effluent nitrogen compounds concentrations. The statistical features of data  
248 sets in group 1 and 2 are listed in **Table 1**.

249 Final data sets in Group 1 to predict effluent nitrogen compounds concentrations resulted in a  
250 total of 1,660 data points, while final data sets in group 2 consisted of 10,956 data points. Same-  
251 day information was classified as “measured” or “calculated” depending on whether data  
252 reflect actual data (40 days) or theoretically calculated data (126 days). Both Group 1 and  
253 Group 2 have 166 groups of data.

### 254 **2.3 Model development and ranking**

255 Random forest (RF) and Support Vector Machine (SVM) are employed for training and testing  
256 of data sets in groups 1 and 2. These algorithms are selected mainly due to their capability of  
257 capturing non-linear relationships in highly complex datasets. RF is an ensemble method which  
258 consists of a user-specified number of decision trees, usually Classification and Regression  
259 Trees, each of which is trained on a sub-sample of the data.<sup>58,59</sup> Decision trees such as CART  
260 are formed by splitting features to make decision so that each split maximizes the reduction of  
261 variance in regression tasks and Gini impurity in classification tasks. SVM models for  
262 regression tasks, also known as Support Vector Regression (SVR) models, are trained by  
263 minimizing a convex cost function that resembles the discrepancy between the actual and  
264 predicted values.<sup>60,61</sup> For comparison purposes, linear regression serves as a baseline to  
265 evaluate the regression performance of these more complex models.

266 The algorithms presented in this paper are developed, trained, and tested using Python 3.6, and  
267 the scikit-learn,<sup>62</sup> TensorFlow<sup>TM</sup>,<sup>63</sup> and KERAS<sup>64</sup> packages. The hyperparameters of the  
268 studied machine learning algorithms were determined using guidelines provided in Asghari et  
269 al.<sup>65</sup> and Wang et al.<sup>66</sup> Optimal values for hyperparameters corresponding to machine learning  
270 algorithms are shown in **Table S10**. All models were trained with 126 theoretically calculated  
271 same-day data and tested with random selections of 20 measured same-day data. This approach  
272 was designed to assess whether models trained on calculated data could achieve comparable

273 performance to those trained on real data. Prior to training neural networks models, all features  
274 were scaled using MinMax function, so they range between 0 and 1. Correlation coefficient  
275 (R), mean absolute error (MAE), and mean squared error (MSE) were used to rank the  
276 performance of the models developed by the machine learning methods.

#### 277 **2.4 Data set size and feature selection analysis**

278 Data size was analyzed in depth and breadth. In depth, because the machine learning algorithm  
279 with the best performance for the data sets was selected to assess whether calculating missing  
280 values through theoretical calculations aids in the performance of the models by varying the  
281 training data set size from 0% (0 same-day data) to 100% (146 same-day data). The R of the  
282 predictions were benchmarked against traditional linear regression. In breadth, because  
283 microbial community abundance percentages were included as input process variables for the  
284 analysis of data sets in group 2. Microbial analysis is labor intensive and time-consuming; thus,  
285 it is vital to determine whether obtaining microbiome data improves predictions. In addition,  
286 the identification of the relative importance (RI) of input process parameters in the prediction  
287 of effluent nitrogen compounds concentrations was performed by the machine learning  
288 algorithm with the best performance. Data sets of Group 1 have nine input process parameters  
289 available, while data sets in Group 2 have over sixty input process parameters. To obtain a  
290 better generalized model, the most relevant parameters were identified. Data sets and machine  
291 learning codes are available online.<sup>1</sup>

---

<sup>1</sup> <https://github.com/vd1371/ML-PNAFMBR>

## 292 3. RESULTS AND DISCUSSION

### 293 3.1 Prediction performance

294 Predicted values for effluent ammonium, nitrite, and nitrate concentrations were compared with  
295 the training, cross validation, and test sets. To simulate the effluent nitrogen compounds  
296 concentrations in the PN/AFMBR, the optimal values for the hyperparameters corresponding  
297 to machine learning algorithms were identified.

298 **Table 2.** Performance of machine learning algorithms on testing data set including microbial community  
299 abundance percentages.

Output process variable	Criteria	LR	RF	SVM
Ammonium	R	0.008	0.894	0.794
	MAE	5.65	2.022	2.540
	MSE	31.96	7.671	13.793
Nitrite	R	0.189	0.722	0.351
	MAE	4.18	0.891	1.389
	MSE	25.84	1.648	5.069
Nitrate	R	0.247	0.802	0.808
	MAE	5.15	1.787	1.857
	MSE	56.14	6.060	6.667
Overall Rank		III	I	II

300

301 Performance evaluations in the test data sets aid to determine the robustness of the predictive  
302 models. **Table 2** displays the R, MAE and MSE values obtained by the three machine learning  
303 methods estimating the difference between effluent and influent ammonium, nitrite and nitrate  
304 concentrations. Ranking of each algorithm based on three evaluation criteria is also provided  
305 in **Table 2**. Based on the results, change in ammonium concentrations are satisfactorily fitted  
306 by random forest with close values for R (0.89), MAE (2.02) and MSE (7.67). For nitrite  
307 concentrations, random forest outperforms the other algorithms by providing the highest R  
308 (0.72), and the lowest MAE (0.89) and MSE (1.65) errors. For the difference in nitrate  
309 concentrations, SVM, presented the highest R of 0.81, while RF presented the lowest MAE of  
310 1.79, and the lowest MSE of 6.06, respectively. Comparison of prediction performance with

311 theoretical-empirical models is not conducted because of the lack of studies modeling change  
312 in nitrogen compounds concentrations in the integrated PN/A FMBR with GAC fluidization.  
313 Indeed, there are some studies that focus on PN/A process without membranes in their  
314 configurations,<sup>17,34,35,67</sup> or MBRs for activated sludge.<sup>31,68</sup> However, some of these studies  
315 focus on mathematical predictions without real data,<sup>17,34</sup> or they lack model validation on test  
316 datasets.<sup>31,35,67,68</sup>

317 Overall, random forest presented the best performance across all nitrogen compounds verifying  
318 its robustness to predict the profile of effluent nitrogen compounds concentrations in the PN/A  
319 FMBR compared to other algorithms. The change in nitrogen compounds concentrations was  
320 accurately predicted within one standard deviation of the data sets in group 1. Unlike SVM  
321 which is trained on the whole training set, elements (trees) of random forest can overlook  
322 mistaken data entries, noises, or anomalies. Moreover, data scaling is not required for the pre-  
323 processing of random forest, unlike SVM. Owing to these facts, training random forest is easier  
324 and more intuitive for engineers. With this flexibility provided by random forest, or ensemble  
325 methods in general, practitioners could, theoretically, train more accurate models on smaller  
326 datasets.

### 327 **3.2 Data size analysis**

328 The comparison of correlation coefficients (R) between predicted concentrations using  
329 different quantities of training data for random forest (RF) and linear regression (LR) is  
330 illustrated in **Figure 2**.  $R^2$  at different sizes of training datasets is included in Supporting  
331 Information, Figure S2. The dataset consists of 166 days in total, including 40 days of measured  
332 same-day data and 126 days of theoretically calculated data for ammonium, nitrite, and nitrate.  
333 Test sets comprised 20 measured same-day data, while training sets varied from 0 to 146 days,

334 with the remaining 20 same-day data in the training set representing approximately 14% of the  
335 data (shown by the blue vertical line).

336 For datasets without microbial community information (**Figure 2A**), the correlation  
337 coefficients stabilized from 50% of the data in the training set onwards. RF consistently  
338 outperformed LR, with correlation ranges of 0.803 to 0.859 for ammonium, 0.665 to 0.772 for  
339 nitrite, and 0.897 to 0.922 for nitrate concentrations, while LR's R values ranged from 0.663  
340 to 0.726 for ammonium, 0.664 to 0.701 for nitrite, and 0.527 to 0.677 for nitrate. The highest  
341 R for ammonium (0.859) and nitrate (0.922) were achieved at 50% of the training data using  
342 RF, while the highest R for nitrite (0.772) was obtained at 80%. Notably, at 40% of the training  
343 set, the R value for LR (0.819) slightly exceeded that of RF (0.815) for ammonium, suggesting  
344 that better tuning of RF hyperparameters could further improve its performance. For nitrite, LR  
345 outperformed RF between 15% and 55% of the training set, while RF showed more stable  
346 performance after 60% of the training data.

347 The performance for nitrate predictions was more consistent with RF, maintaining correlation  
348 values between 0.860 and 0.922, while LR exhibited decreasing R values with increasing data  
349 size, dropping from 0.527 to 0.882. The lowest R values for both algorithms occurred when  
350 the training set was limited to the 20 measured same-day data. However, training sets  
351 containing 50% to 60% of the data (73 to 88 same-day data) reached acceptable R values using  
352 RF, highlighting the importance of calculated values for missing data.

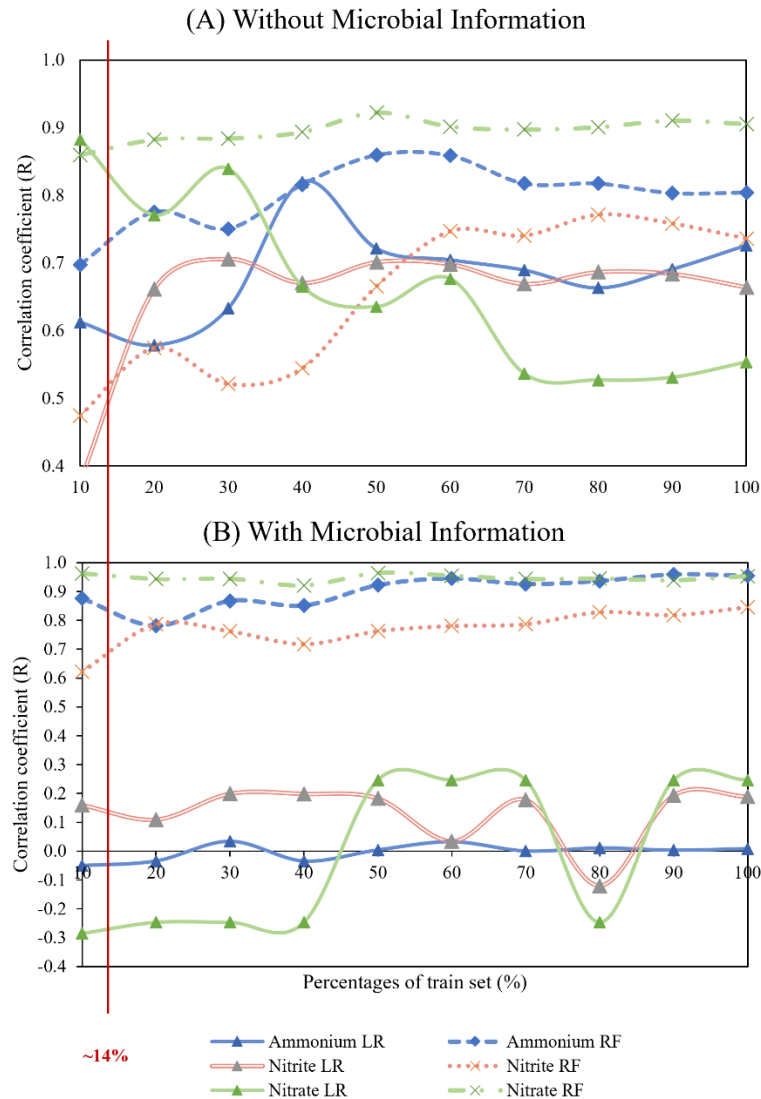
353 When microbial community information was included (**Figure 2B**), the model performance  
354 improved significantly. For ammonium, the R value increased from 0.894 to 0.938, mean  
355 absolute error (MAE) decreased from 2.022 to 1.572, and mean squared error (MSE) dropped  
356 from 7.671 to 4.508, as presented in **Table S11**. For nitrite, R improved from 0.722 to 0.851,  
357 with MAE decreasing from 0.891 to 0.651, and MSE from 1.648 to 1.346. Similarly, nitrate



358 predictions improved, with R increasing from 0.802 to 0.942, MAE decreasing from 1.787 to  
359 0.917, and MSE from 6.060 to 1.915. This demonstrates that including microbial community  
360 data enhances model accuracy by providing more comprehensive system information.

361 With the additional microbial data, RF continued to show strong performance, with consistent  
362 R values ranging from 0.781 to 0.959 for ammonium, 0.621 to 0.846 for nitrite, and 0.920 to  
363 0.964 for nitrate, compared to LR. The highest R values were obtained at 90%, 100%, and 50%  
364 of the data for ammonium, nitrite, and nitrate, respectively. For ammonium and nitrate, R  
365 values greater than 0.90 were achieved with only 50% and 10% of the training data,  
366 respectively, while nitrite predictions required larger data sets to reach an R value above 0.80.

367 The hybrid theoretical-machine learning approach leverages stoichiometry and kinetic models  
368 to fill in missing data for empirical analysis, offering better predictions for nitrogen  
369 concentrations. Besides, the comparison with LR indicates that ML-Theoretical Model deliver  
370 a better performance even under constraints of limited data, and the predictive performance  
371 improves faster than LR when the data size increases. While traditional theoretical models may  
372 overlook complex microbial interactions, integrating microbial community data allowed this  
373 approach to identify empirical patterns and interactions that theoretical calculations alone  
374 might miss. This combination of methods leads to a more comprehensive understanding of  
375 system dynamics, as demonstrated by the improved predictions over the First Principal Model.  
376 Detailed comparisons between the ML-theoretical approach and the First Principal Model,  
377 including the effects on prediction outcomes, are provided in **Supporting Method S3 and**  
378 **Tables S2-S7.**



379

380 **Figure 2.** Correlation coefficients at different sizes of training data sets as predicted by linear regression (LR)  
 381 and random forest (RF). Red vertical line represents the percentage of test data sets (~14%). (A) datasets  
 382 without microbial community information; (B) datasets include microbial community information.

383

### 384 3.3 Identification of Key Environmental Factors

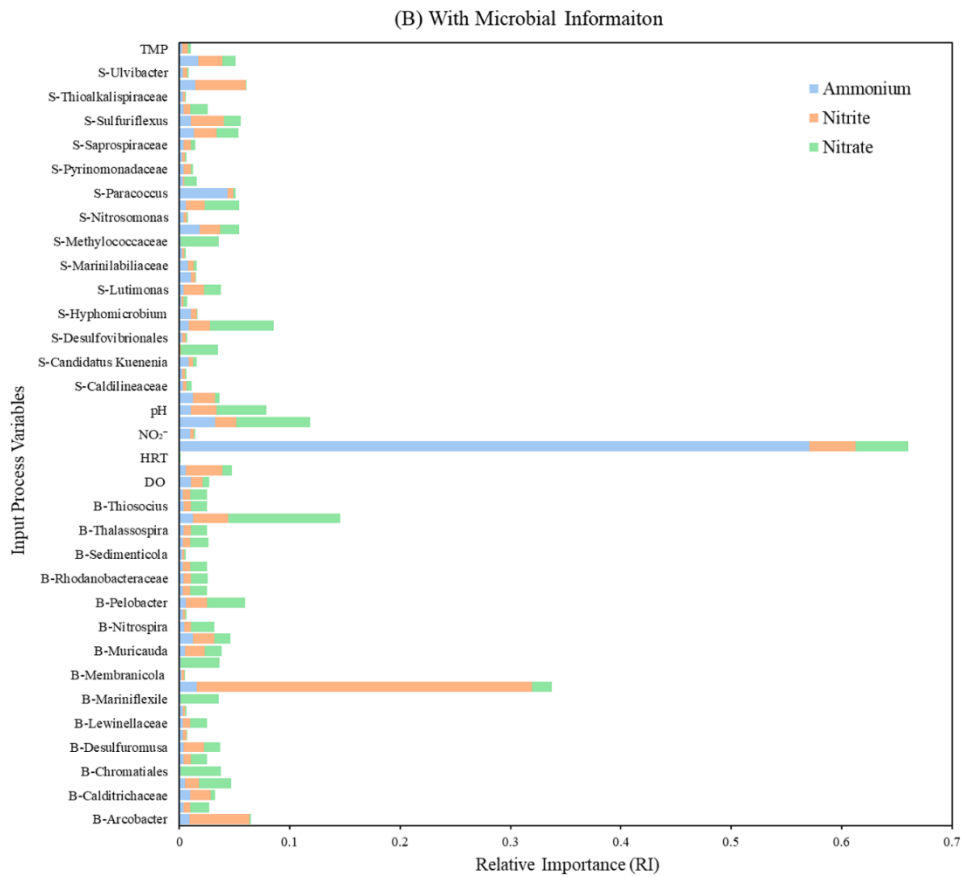
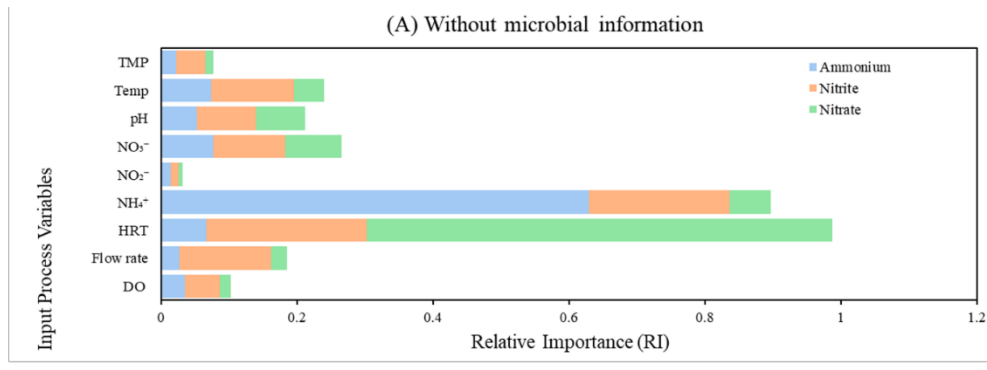
385 Variations from influent to effluent nitrogen compounds concentration depends on several  
 386 parameters. For the change in ammonium concentrations in Group 1 (**Figure 3A**), pollution  
 387 load (ammonium with RI of 0.630 and nitrate concentration in the influent with RI of 0.077),  
 388 and temperature (RI = 0.074) were the dominant parameters influencing the prediction values.

389 These three parameters were also influential when considering microbial community  
390 abundances (Group 2) with RI of 0.570 for ammonium concentration in the influent, of 0.032  
391 for nitrate concentration in the influent, and of 0.017 for temperature.

392 Temperature, known to be critical for nitrification, is consistent with previous findings <sup>40</sup>,  
393 especially in relation to the optimal growth rate of *Nitrosomonas* (25-35°C). Feature selection  
394 from random forest algorithms also revealed that nitrate concentration in the influent  
395 significantly impacts ammonium concentration predictions. Feature selection from random  
396 forest algorithms also revealed that nitrate concentration in the influent significantly impacts  
397 ammonium concentration predictions.

398 For nitrite concentrations, HRT (RI = 0.236), influent ammonium (RI = 0.207), flow rate (RI  
399 = 0.135), temperature (RI = 0.122), and influent nitrate (RI = 0.106) were identified as key  
400 factors in group 1. These environmental factors were less dominant when microbial data was  
401 incorporated as shown in group 2 where influent ammonium concentration (RI = 0.042), flow  
402 rate (RI = 0.032), and temperature (RI = 0.021) remained as influential parameters (**Figure 3B**).  
403 The relative importance of HRT decreased because microbes – especially in biofilm – depend  
404 in solids retention time rather than HRT

405 For nitrate concentrations, HRT (RI = 0.685), nitrate concentration in the influent (RI = 0.083),  
406 and pH (RI = 0.072) emerged as the most critical parameters in Group 1. In group 2, nitrate  
407 concentration in the influent (RI = 0.067), ammonium concentration in the influent (RI =  
408 0.047), and pH (RI = 0.045) are the key operational parameters identified (**Figure 3B**). Similar  
409 to nitrite data sets, the relative importance of HRT decreased because microbes are not  
410 controlled by HRT but solids retention time. Microbiomes removing nitrate presented higher  
411 influence in nitrate concentrations than microbiomes producing nitrate.



412

413  
414

**Figure 3.** Relative Importance of influential parameters for effluent concentrations: Ammonium, Nitrite concentrations, Nitrate (A) without microbial information (B) with microbial information.

### 415 3.4 Discussion of Microbial Interactions

416 In Group 2 (**Figure 3B**), which considered microbial community data, specific bacteria played  
417 critical roles in nitrogen transformations. For ammonium concentrations, *Paracoccus* (RI =  
418 0.043), *Muricauda* (RI = 0.018), *Thiosocius* (RI = 0.015), *Sedimenticola* (RI = 0.013),  
419 *Arcobacter* (RI = 0.012), and *Hyphomicrobium* (RI = 0.011) in suspension sludge, and  
420 *Marinilabiliaceae* (RI = 0.016), *Thioalkalspiraceae* (RI = 0.013) and *Nitrosomonas* (RI =  
421 0.013) in the biofilm were also identified as influential factors in the change of ammonium  
422 concentrations in group 2. Given that ammonium is removed during the first step of nitrification  
423 as performed by ammonium oxidizing bacteria (i.e., *Nitrosomonas*), influent ammonium  
424 concentration and temperature are expected as key influential parameters for change in  
425 ammonium concentrations predictions. *Hyphomicrobium*<sup>70</sup>, *Sedimenticola*<sup>71</sup>, *Thiosocius*<sup>71</sup>,  
426 and *Arcobacter*<sup>72</sup> in suspension sludge; and *Marinilabiliaceae*<sup>73</sup> and *Thioalkalspiraceae*<sup>74</sup> in  
427 the biofilm of GAC reduce nitrate while removing or producing low quantities of ammonium.  
428 Positive correlations were previously found through metagenomic sequencing analyses  
429 between *Hyphomicrobium* and *Paracoccus*, and *Sedimenticola* and *Muricauda*.<sup>49</sup> *Paracoccus*<sup>75</sup>  
430 and *Muricauda* remove organics while producing low quantities of ammonium concentrations.  
431 *Nitrosomonas* (ammonium oxidizing bacteria) with nitrite production, and *Nitrospira* (NOB)  
432 and *Candidatus Kuenenia* (Anammox bacteria) with nitrite removal were expected to lead as  
433 key microbiomes in the change of nitrite concentrations. However, random forest algorithm  
434 suggests that partial denitrifiers have higher influence than *Nitrosomonas* (RI = 0.019 in  
435 biofilm), *Nitrospira* (RI = 0.017 in suspension sludge) and *Anammox* bacteria (RI = 0.012 in  
436 biofilm) in this PN/AFMBR. This could be due to several reasons. First, partial denitrifiers may  
437 be more adaptable to the varying environmental conditions within the reactor, allowing them  
438 to maintain higher activity levels<sup>76</sup>. Second, the metabolic flexibility of partial denitrifiers

439 might enable them to utilize a broader range of substrates including both inorganic and organic  
440 carbon, thus exerting a more significant influence on the overall nitrogen removal process.  
441 Lastly, interactions between partial denitrifiers and other microbial communities could enhance  
442 their relative importance, possibly through synergistic relationships that boost their  
443 effectiveness in nitrogen compound transformation. Further investigation is needed to elucidate  
444 these interactions and their impact on the system's performance.

445 For nitrite concentration, *Marinilabiliaceae* (RI = 0.304) and *Arcobacter* (RI = 0.055) in the  
446 biofilm, and *Thiosocius* (RI = 0.045) in the suspensions sludge were the dominant factors.  
447 *Marinilabiliaceae*<sup>73</sup> and *Arcobacter*<sup>72</sup> have been identified as heterotrophic partial denitrifiers  
448 with nitrate reduction to nitrite. Positive correlations between *Thiosocius* and *Arcobacter* in  
449 suspension sludge, and with *Marinilabiliaceae* through intermediate species (i.e.,  
450 *Rhodanobacteraceae*) in biofilm were identified through previous metagenomic sequencing  
451 analysis.<sup>49</sup> *Sulfuriflexus*<sup>77</sup> and *Sedimenticola*<sup>71</sup> have similar physiological functions to  
452 *Thiosocius*<sup>78</sup> and have been identified as autotrophic partial denitrifiers with nitrate reduction.  
453 *Sulfuriflexus* (RI = 0.030) and *Sedimenticola* (RI = 0.020) in the suspension sludge were also  
454 key influential microbiome. Nitrite has not been found to be consumed nor produced by  
455 *Thioalkalispiraceae*<sup>74</sup> (RI = 0.031) and *Pelobacter*<sup>77</sup> (RI = 0.019); yet random forest method  
456 found these species in the biofilm to be influential in the change of nitrite concentrations.

457 In nitrate predictions, *Thioalkalispiraceae* (RI = 0.102) and *Chromatiales* (RI = 0.037) in the  
458 biofilm, and *Chromatiales* (RI = 0.034) in the suspension sludge are autotrophic denitrifiers  
459 removing nitrate.<sup>74,79</sup> *Chromatiales* demonstrated positive correlation to *Methylococcaceae* as  
460 investigated by metagenomics sequencing analysis.<sup>80</sup> *Methylococcaceae* present in biofilm (RI  
461 = 0.035) and suspension sludge (RI = 0.035) reduce nitrate to nitric oxide through methanol  
462 oxidation.<sup>81,82</sup> *Desulfuromusa* (RI = 0.058) in suspension sludge, and *Pelobacter* (RI = 0.034)  
463 in biofilm are key microbiome consuming nitrate through heterotrophic dissimilatory nitrate

464 reduction to ammonium.<sup>77,83</sup> *Nitrospira* in suspension sludge (RI = 0.031), *Candidatus*  
465 *Kuenenia* (RI = 0.029) and *Nitrospira* in biofilm (RI = 0.021) are nitrate producing bacteria  
466 which influenced the change in nitrate concentration as expected.<sup>5</sup> This also aligns with the  
467 understanding that maintaining stable suppression of nitrite-oxidizing bacteria (NOB) is crucial  
468 for achieving optimal PN/A performance.<sup>84</sup> *Mariniflexile*<sup>85</sup> in biofilm (RI = 0.035) does not  
469 consume nor produce nitrate but presented positive correlation to *Thioalkalispiraceae* in  
470 biofilm which might indicate synergistic relationships among them.

471 Previous PN/AFMBR model groups microbiomes depending on their physiological  
472 functions.<sup>49</sup> For example, heterotrophic partial denitrifiers with nitrate reduction to nitrite  
473 within the PN/AFMBR included *Arcobacter*, *Ulvibacter*, *Lutimonas*, and *Marinilabiliaceae*.  
474 This hybrid approach including machine learning algorithms provided in-depth analysis by  
475 identifying that among all heterotrophic partial denitrifiers, *Marinilabiliaceae* and *Arcobacter*  
476 in the biofilm are the key influential microbiomes affecting the change of nitrite concentrations  
477 in saline sewage. Microbiomes identified at the genus level obtained by 16S rRNA gene  
478 sequencing are not handled by traditional biological models. Yet, data sets for statistical-  
479 empirical approach would contain missing data without fixed theoretical calculations as  
480 provided in traditional models. Thus, the integration of both approaches overcomes each  
481 other's shortcomings and identifies hidden microbial interactions.

### 482 **3.5 Environmental Implications**

483 This study proposed a hybrid theoretical-machine learning approach framework, highlighting  
484 the feasibility of leveraging limited microbial data for accurate predictions of nitrogen removal  
485 performance. By integrating microbial community data, the framework provides actionable  
486 insights into the key environmental parameters and microbial interactions that influence  
487 nitrogen removal efficiency. Moreover, this research underscores the importance of improving

488 microbial data collection methods and integrating them into predictive models. Enhanced  
489 microbial monitoring in full-scale systems could facilitate more precise control over nitrogen  
490 removal processes, leading to better environmental outcomes. By bridging the gap between  
491 pilot-scale experimentation and real-world applications, this study paves the way for the  
492 development of advanced decision-making tools tailored to biological nitrogen removal  
493 systems.

### 494 **3.6 Limitations & future works**

495 This hybrid method of theoretical calculations and machine learning algorithms has certain  
496 limitations. First, theoretical models depend on known microbiomes and steady-state  
497 assumptions, which can miss dynamic changes in microbial populations. Despite advances in  
498 meta-omics, about 27-60% of the microbiomes in the PN/A FMBR samples were “unclassified.”  
499 As new bacteria are discovered, they can be incorporated into models like the TEE, improving  
500 predictions by updating missing data. Steady-state models also struggle with daily fluctuations  
501 in influent characteristics, requiring random values to fill gaps for unmeasured periods. Besides,  
502 machine learning algorithms need large datasets to train effectively. In this study, 73-88 same-  
503 day measurements were necessary to achieve high correlation predictions for nitrogen  
504 compounds, but collecting and analyzing microbial DNA from such large samples can be  
505 expensive and time-consuming. However, the hybrid approach mitigates individual limitations  
506 by leveraging the strengths of both methods, resulting in acceptable prediction errors using the  
507 random forest algorithm.

508 Given that random forest algorithms presented well fitted predictions of changes in nitrogen  
509 compounds concentrations of a PN/A FMBR, this hybrid theoretical-machine learning  
510 approach could be tested for other biological wastewater and sludge treatment processes.  
511 Furthermore, hidden microbiomes relationships could be further investigated. Future research



512 should explore hidden microbial relationships, particularly interactions between key species  
513 like *Thiosocius* and *Marinilabiliaceae* in the biofilm, and further investigate how autotrophic  
514 denitrifiers like *Thioalkalispiraceae* interact with heterotrophic reducers like *Pelobacter*.  
515 Additionally, the study identified 11 key genera in both suspension sludge and biofilm that  
516 influence nitrogen concentrations, suggesting the need for further exploration of their  
517 interactions with operational parameters. Future research should focus on validating this  
518 methodology with additional datasets from full-scale systems to ensure its robustness and  
519 generalizability. This validation will help refine the model's predictive accuracy and further  
520 optimize operational strategies. As full-scale systems often encounter similar challenges, such  
521 as data gaps and fluctuating influent characteristics, this study provides a promising approach  
522 to addressing these limitations.

#### 523 **4. CONCLUSIONS**

524 Traditional sewage treatment processes often face difficulties due to incomplete or invalid data,  
525 which hampers effective optimization and scaling. By utilizing a hybrid theoretical-machine  
526 learning approach, this study successfully bridges these gaps by calculating missing values and  
527 revealing hidden microbial interactions that traditional methods overlook. The main  
528 conclusions are as follows:

- 529 1) Random forest consistently outperformed other algorithms, showing robust prediction  
530 accuracy for ammonium, nitrite, and nitrate concentrations, with correlation  
531 coefficients of 0.89, 0.72, and 0.80, respectively.
- 532 2) Including microbial community abundance data significantly improved model  
533 performance, revealing critical microbial interactions, particularly for nitrogen  
534 transformation and partial denitrification.

535 3) Key parameters such as temperature, pH, and microbial communities (e.g.,  
536 heterotrophic denitrifiers like *Paracoccus*, *Marinilabiliaceae* and autotrophic  
537 denitrifiers, such as *Thioalkalispiraceae*) were identified as essential for efficient  
538 PN/AFMBR performance.

## 539 5. ACKNOWLEDGEMENTS

540 The authors thank the Hong Kong Research Grants Council- University Grants Committee  
541 (Grant No. 15252916 and UGC/GEN/456/08), and the Research Institute for Sustainable Urban  
542 Development (RISUD) for their financial support. Declarations of interest: none.

## 543 6. SUPPORTING INFORMATION

544 The “Supporting Information” file contains the assumptions for missing values of nitrogen  
545 compounds concentrations in the effluent, explanation of dates and frequency of microbial  
546 community analysis sampling, prediction of effluent parameters based on mathematical  
547 models, statistical description of microbial community abundance percentages data, optimal  
548 values of hyperparameters for machine learning algorithms and performance of machine  
549 learning algorithms on test set including microbial community abundance percentages.

## 550 7. REFERENCES

- 551 (1) McCarty, P. L. What is the Best Biological Process for Nitrogen Removal: When and  
552 Why? *Environ. Sci Technol* **2018**, *52*, 3835–3841.
- 553 (2) Chen, R.; Ji, J.; Chen, Y.; Takemura, Y.; Liu, Y.; Kubota, K.; Ma, H.; Li, Y.-Y. Successful  
554 Operation Performance and Syntrophic Micro-Granule in Partial Nitritation and  
555 Anammox Reactor Treating Low-Strength Ammonia Wastewater. *Water Res.* **2019**, *155*,  
556 288–299. <https://doi.org/10.1016/j.watres.2019.02.041>.
- 557 (3) Huang, X.; Sun, K.; Wei, Q.; Urata, K.; Yamashita, Y.; Hong, N.; Hama, T.; Kawagoshi,  
558 Y. One-Stage Partial Nitritation and Anammox in Membrane Bioreactor. *Env. Sci Pollut*  
559 *Res* **2016**, *23*, 11149–11162.
- 560 (4) Huang, X.; Mi, W.; Chan, Y. H.; Singh, S.; Zhuang, H.; Leu, S. Y.; Li, X.; Li, X.; Lee, P.  
561 H. C-N-S Synergy in a Pilot-Scale Mainstream Anammox Fluidized-Bed Membrane

- 562 Bioreactor for Treating Chemically Enhanced Primary Treatment Saline Sewage. *Water*  
563 *Res.* **2023**, 229 (1), 119475.
- 564 (5) Jetten, M. S. M.; Horn, S. J.; Loosdrecht, M. C. M. Towards a More Sustainable  
565 Municipal Wastewater Treatment System. *Water Sci Technol* **1997**, 35, 171–180.
- 566 (6) Mulder, A.; Graaf, A. A.; Robertson, L. A.; Kuenen, J. G. Anaerobic Ammonium  
567 Oxidation Discovered in Nitrifying Fluidized Bed Reactor. *Microbiol Ecol* **1995**, 16, 177–  
568 184.
- 569 (7) Strous, M.; Heijnen, J. J.; Kuenen, J. G.; Jetten, M. S. M. The Sequencing Batch Reactor  
570 as a Powerful Tool for the Study of Slowly Growing Anaerobic Ammonium-Oxidizing  
571 Microorganisms. *Appl. Microbiol. Biotechnol.* **1998**, 50 (5), 589–596.  
572 <https://doi.org/10.1007/s002530051340>.
- 573 (8) Lee, P.; Kwak, W.; Bae, J.; Mccarty, P. L. The Effect of SRT on Nitrate Formation during  
574 Autotrophic Nitrogen Removal of Anaerobically Treated Wastewater. *Water Sci Technol*  
575 **2013**, 68, 1751–1756.
- 576 (9) Abdelmelek, S. B.; Greaves, J.; Ishida, K. P.; Cooper, W. J.; Song, W. Removal of  
577 Pharmaceutical and Personal Care Products from Reverse Osmosis Retentate Using  
578 Advanced Oxidation Processes. *Environ. Sci. Technol.* **2011**, 45 (8), 3665–3671.  
579 <https://doi.org/10.1021/es104287n>.
- 580 (10) Katsou, E.; Malamis, S.; Loizidou, M. Performance of a Membrane Bioreactor Used for  
581 the Treatment of Wastewater Contaminated with Heavy Metals. *Bioresour Technol* **2011**,  
582 102, 4325–4332.
- 583 (11) Wu, B.; Li, Y.; Lim, W.; Lin, S.; Guo, Q.; Fane, A. G. Single-Stage versus Two-Stage  
584 Anaerobic Fluidized Bed Bioreactors in Treating Municipal Wastewater : Performance,  
585 Foulant Characteristics , and Microbial Community. *Chemosphere* **2017**, 171, 158–167.
- 586 (12) Chen, F.; Qian, Y.; Cheng, H.; Shen, J.; Qin, Y.; Li, Y.-Y. Recent Developments in  
587 Anammox-Based Membrane Bioreactors: A Review. *Sci. Total Environ.* **2023**, 857,  
588 159539. <https://doi.org/10.1016/j.scitotenv.2022.159539>.
- 589 (13) Li, J.; Li, J.; Peng, Y.; Wang, S.; Zhang, L.; Yang, S.; Li, S. Insight into the Impacts of  
590 Organics on Anammox and Their Potential Linking to System Performance of Sewage  
591 Partial Nitrification-Anammox (PN/A): A Critical Review. *Bioresour. Technol.* **2020**,  
592 300, 122655. <https://doi.org/10.1016/j.biortech.2019.122655>.
- 593 (14) Wang, Z.; Zheng, M.; Duan, H.; Yuan, Z.; Hu, S. A 20-Year Journey of Partial Nitritation  
594 and Anammox (PN/A): From Sidestream toward Mainstream. *Environ. Sci. Technol.*  
595 **2022**, 56 (12), 7522–7531. <https://doi.org/10.1021/acs.est.1c06107>.
- 596 (15) Zhang, X.; Zhang, X.; Chen, J.; Wu, P.; Yang, Z.; Zhou, L.; Zhu, Z.; Wu, Z.; Zhang, K.;  
597 Wang, Y.; Ruth, G. A Critical Review of Improving Mainstream Anammox Systems:  
598 Based on Macroscopic Process Regulation and Microscopic Enhancement Mechanisms.  
599 *Environ. Res.* **2023**, 236, 116770. <https://doi.org/10.1016/j.envres.2023.116770>.
- 600 (16) Harrou, F.; Dairi, A.; Sun, Y.; Senouci, M. Statistical Monitoring of a Wastewater  
601 Treatment Plant: A Case Study. *J Env. Manage* **2018**, 223, 807–814.
- 602 (17) Hubaux, N.; Wells, G.; Morgenroth, E. Impact of Coexistence of Flocs and Biofilm on  
603 Performance of Combined Nitritation-Anammox Granular Sludge Reactors. *Water Res*  
604 **2015**, 68, 127–139.

- 605 (18) Rodriguez-Sanchez, A.; Muñoz-Palazon, B.; Hurtado-Martinez, M.; Maza-Marquez, P.;  
606 Gonzalez-Lopez, J.; Vahala, R.; Gonzalez-Martinez, A. Microbial Ecology Dynamics of  
607 a Partial Nitrification Bioreactor with Polar Arctic Circle Activated Sludge Operating at  
608 Low Temperature. *Chemosphere* **2019**, *225*, 73–82.
- 609 (19) Al-Hazmi, H. E.; Lu, X.; Grubba, D.; Majtacz, J.; Badawi, M.; Mąkinia, J. Sustainable  
610 Nitrogen Removal in Anammox-Mediated Systems: Microbial Metabolic Pathways,  
611 Operational Conditions and Mathematical Modelling. *Sci. Total Environ.* **2023**, *868*,  
612 161633. <https://doi.org/10.1016/j.scitotenv.2023.161633>.
- 613 (20) Sharifshourjeh, M.; Kowal, P.; Xi, L.; Xie, L.; Drewnowski, J. Development of Strategies  
614 for AOB and NOB Competition Supported by Mathematical Modeling in Terms of  
615 Successful Deammonification Implementation for Energy-Efficient WWTPs. *Processes*  
616 **2021**, *9*, 562. <https://doi.org/10.3390/pr9030562>.
- 617 (21) Kosgey, K.; Zungu, P. V.; Kumari, S.; Bux, F. Critical Review of Process Control  
618 Strategies in Anammox-Mediated Nitrogen Removal Systems. *J. Environ. Chem. Eng.*  
619 **2022**, *10* (4), 108068. <https://doi.org/10.1016/j.jece.2022.108068>.
- 620 (47) Lamrini, B.; Lakhal, E.-K.; Le Lann, M.-V.; Wehenkel, L. Data Validation and Missing  
621 Data Reconstruction Using Self-Organizing Map for Water Treatment. *Neural Comput.*  
622 *Appl.* **2011**, *20* (4), 575–588. <https://doi.org/10.1007/s00521-011-0526-5>.
- 623 (23) Vannecke, T. P. W.; Bernet, N.; Winkler, M. K. H.; Santa-Catalina, G.; Steyer, J. P.;  
624 Volcke, E. I. P. Influence of Process Dynamics on the Microbial Diversity in a Nitrifying  
625 Biofilm Reactor: Correlation Analysis and Simulation Study. *Biotechnol Bioeng* **2016**,  
626 *113*, 1962–1974.
- 627 (24) Gujer, W.; Henze, M.; Mino, T.; Loosdrecht, M. Activated sludge model no.3. *Water Sci*  
628 *Technol* **1999**, *39*, 183–193.
- 629 (25) Henze, M.; Grady Jr, L.; Gujer, W.; Marais, G.; Matsuo, T. Activated Sludge Model No  
630 1. *Wat Sci Technol* **1987**, *29*.
- 631 (26) Batstone, D. J.; Keller, J.; Angelidaki, I.; Kalyuzhnyi, S. V.; Pavlostathis, S. G.  
632 *Anaerobic digestion model no.1 (ADM1)*; Water Science and technology, 2002; Vol.  
633 45(10).
- 634 (27) Reichert, P. AQUASIM 2.0 - User Manual. In *Computer Program for the Identification*  
635 *and Simulation of Aquatic Systems*; 1998.
- 636 (28) Water and Wastewater Treatment Modeling and Simulation Software| Hydromantis  
637 VERSION 8.0, 2019. <https://www.hydromantis.com/> (accessed 2023-12-28).
- 638 (29) *EnviroSim – Wastewater Modeling Software*. <https://envirosim.com/> (accessed 2023-12-  
639 28).
- 640 (30) Cho, J. W.; Ahn, K. H.; Lee, Y. H.; Lim, B. R.; Kim, J. Y. Investigation of Biological and  
641 Fouling Characteristics of Submerged Membrane Bioreactor Process for Wastewater  
642 Treatment by Model Sensitivity Analysis. *Water Sci Technol* **2004**, *49*, 245–254.
- 643 (31) Sarioglu, M.; Insel, G.; Artan, N.; Orhon, D. Stoichiometric and Kinetic Evaluation of  
644 Simultaneous Nitrification and Denitrification in a Membrane Bioreactor at Steady State.  
645 *J Chem Technol Biotechnol* **2011**, *86*, 798–811.
- 646 (32) Mehrdad, M.; Park, H.; Chandran, K.; Ramalingam, K.; Fillos, J. Biofilm Population  
647 Diversity and Distribution in an Anammox MBBR Pilot: Molecular Analysis and

- 648 Mathematical Modeling; Water Environment Federation, 2016; Vol. 2016/11.  
649 <https://doi.org/10.2175/193864716819706707>.
- 650 (33) Pérez, J.; Lotti, T.; Kleerebezem, R.; Picioreanu, C.; Loosdrecht, M. C. M. Outcompeting  
651 Nitrite-Oxidizing Bacteria in Single-Stage Nitrogen Removal in Sewage Treatment  
652 Plants: A Model-Based Study. *Water Res* **2014**, *66*, 208–218.
- 653 (34) Rittmann, B. E.; Boltz, J. P.; Brockmann, D.; Daigger, G. T.; Morgenroth, E.; Sørensen,  
654 K. H.; Takács, I.; Loosdrecht, M.; Vanrolleghem, P. A. A Framework for Good Biofilm  
655 Reactor Modeling Practice (GBRMP). *Water Sci. Technol* **2018**, *77*, 1149–1164.
- 656 (35) Volcke, E. I. P.; Picioreanu, C.; Baets, B.; Loosdrecht, M. C. M. The Granule Size  
657 Distribution in an Anammox-Based Granular Sludge Reactor Affects the Conversion-  
658 Implications for Modeling. *Biotechnol Bioeng* **2012**, *109*, 1629–1636.
- 659 (36) Kim, M.; Kim, Y.; Kim, H.; Piao, W.; Kim, C. Evaluation of the K-Nearest Neighbor  
660 Method for Forecasting the Influent Characteristics of Wastewater Treatment Plant. *Front*  
661 *Env. Sci Eng* **2016**, *10*, 299–310.
- 662 (37) Newhart, K. B.; Holloway, R. W.; Hering, A. S.; Cath, T. Y. Data-Driven Performance  
663 Analyses of Wastewater Treatment Plants: A Review. *Water Res.* **2019**, *157*, 498–513.  
664 <https://doi.org/10.1016/j.watres.2019.03.030>.
- 665 (38) Park, M.; Anumol, T.; Snyder, S. A. Modeling Approaches to Predict Removal of Trace  
666 Organic Compounds by Ozone Oxidation in Potable Reuse Applications. *Env. Sci Water*  
667 *Res Technol* **2015**, *1*, 699–708.
- 668 (39) Qiao, J.; Hu, Z.; Li, W. Soft Measurement Modeling Based on Chaos Theory for  
669 Biochemical Oxygen Demand. *Water* **2016**, *8(12)* (581).
- 670 (40) Verma, A.; Wei, X.; Kusiak, A. Predicting the Total Suspended Solids in Wastewater: A  
671 Data-Mining Approach. *Eng Appl Artif Intell* **2013**, *26*, 1366–1372.
- 672 (41) Vega, P. T.; Jaramillo-Moran, M. A. Obtaining Key Parameters and Working Conditions  
673 of Wastewater Biological Nutrient Removal by Means of Artificial Intelligence Tools.  
674 *Water* **2018**, *10*.
- 675 (42) Awolusi, O. O.; Nasr, M.; Kumari, S.; Bux, F. Artificial Intelligence for the Evaluation  
676 of Operational Parameters Influencing Nitrification and Nitrifiers in an Activated Sludge  
677 Process. *Microb Ecol* **2016**, *72*, 49–63.
- 678 (43) Zhao, L.; Wu, Q. Review and Expectation of Artificial Intelligent System for Wastewater  
679 Treatment. *Appl Mech Mater* **2013**, *422*, 237–241.
- 680 (44) Wang, D.; Thunéll, S.; Lindberg, U.; Jiang, L.; Trygg, J.; Tysklind, M. Towards Better  
681 Process Management in Wastewater Treatment Plants: Process Analytics Based on SHAP  
682 Values for Tree-Based Machine Learning Methods. *J. Environ. Manage.* **2022**, *301*,  
683 113941. <https://doi.org/10.1016/j.jenvman.2021.113941>.
- 684 (45) Xu, Y.; Zeng, X.; Bernard, S.; He, Z. Data-Driven Prediction of Neutralizer pH and Valve  
685 Position towards Precise Control of Chemical Dosage in a Wastewater Treatment Plant.  
686 *J. Clean. Prod.* **2022**, *348*, 131360. <https://doi.org/10.1016/j.jclepro.2022.131360>.
- 687 (46) Olsson, G. Instrumentation, Control and Automation in the Water Industry – State-of-the-  
688 Art and New Challenges. *Water Sci. Technol.* **2006**, *53* (4–5), 1–16.  
689 <https://doi.org/10.2166/wst.2006.097>.

- 690 (47) Lamrini, B.; Lakhal, E.-K.; Le Lann, M.-V.; Wehenkel, L. Data Validation and Missing  
691 Data Reconstruction Using Self-Organizing Map for Water Treatment. *Neural Comput.*  
692 *Appl.* **2011**, *20* (4), 575–588. <https://doi.org/10.1007/s00521-011-0526-5>.
- 693 (48) Duarte, M. S.; Martins, G.; Oliveira, P.; Fernandes, B.; Ferreira, E. C.; Alves, M. M.;  
694 Lopes, F.; Pereira, M. A.; Novais, P. A Review of Computational Modeling in  
695 Wastewater Treatment Processes. *ACS EST Water* **2024**, *4* (3), 784–804.  
696 <https://doi.org/10.1021/acsestwater.3c00117>.
- 697 (49) Eaton, A. D.; Clesceri, L. S.; Franson, M. A. H.; Rice, E. W.; Greenberg, A. E. Standard  
698 Methods for the Examination of Water and Wastewater, 2005.
- 699 (50) Kanehisa, M.; Sato, Y.; Morishima, K. BlastKOALA and GhostKOALA: KEGG Tools  
700 for Functional Characterization of Genome and Metagenome Sequences. *J. Mol. Biol.*  
701 **2016**, *428*, 4, 726–731.
- 702 (51) Torregrossa, D.; Leopold, U.; Hernández-Sancho, F.; Hansen, J. Machine Learning for  
703 Energy Cost Modelling in Wastewater Treatment Plants. *J. Env. Manage* **2018**, *223*, 1061–  
704 1067.
- 705 (52) Asadi, A.; Verma, A.; Yang, K.; Mejabi, B. Wastewater Treatment Aeration Process  
706 Optimization: A Data Mining Approach. *J. Env. Manage* **2017**, *203*, 630–639.
- 707 (53) Drainage Services Department. Data Collected Directly from the Government  
708 Department in Hong Kong S.A.R, 2017.
- 709 (54) Newhart, K. B.; Holloway, R. W.; Hering, A. S.; Cath, T. Y. Data-Driven Performance  
710 Analyses of Wastewater Treatment Plants: A Review. *Water Res.* **2019**, *157*, 498–513.  
711 <https://doi.org/10.1016/j.watres.2019.03.030>.
- 712 (55) Chaoui, A.; Rebija, K.; Chkaiti, K.; Laaouan, M.; Bourziza, R.; Sebari, K.; Elkhoumsi,  
713 W. Applications of Missing Data Imputation Methods in Wastewater Treatment Plants A  
714 Systematic Literature Review Using the Kitchenham Method. *Int. J. Adv. Comput. Sci.*  
715 *Appl.* **2023**, *14*, 461. <https://doi.org/10.14569/IJACSA.2023.0141049>.
- 716 (56) Tieleman, T.; Hinton, G. Lecture 6.5-Rmsprop: Divide the Gradient by a Running  
717 Average of Its Recent Magnitude. *COURSERA Neural Netw. Mach. Learn.* **2012**, *4*, 2.
- 718 (57) Kingma, D. P.; Ba, J. L. Adam: A Method for Stochastic Optimization. In *3rd*  
719 *International Conference on Learning Representations, ICLR 2015 - Conference Track*  
720 *Proceedings*; 2015; pp 1–15.
- 721 (58) Breiman, L. *Random For.* **2001**, *45*, 5–32.
- 722 (59) Breiman, L.; Friedman, J.; Olshen, R. A.; Stone, C. J. *Classification and Regression*  
723 *Trees*; Chapman and Hall/CRC, 1984. <https://doi.org/10.1201/9781315139470>.
- 724 (60) Drucker, H.; Surges, C. J. C.; Kaufman, L.; Smola, A.; Vapnik, V. Support Vector  
725 Regression Machines. *Adv Neural Inf Process Syst* **1997**, *1*, 155–161.
- 726 (61) Cristianini, N.; Ricci, E. Support Vector Machines, 2008.
- 727 (62) Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel,  
728 M.; Prettenhofer, P.; Weiss, R.; Dubourg, V.; Vanderplas, J.; Passos, A.; Cournapeau, D.;  
729 Brucher, M.; M., P.; Duchesnay, E. Scikit-Learn: Machine Learning in Python. *J. Mach*  
730 *Learn Res* **2011**, *12*, 2825–2830.
- 731 (63) Abadi, M.; Agarwal, A.; Barham, P.; Brevdo, E.; Chen, Z.; Citro, C.; Corrado, G. S.;  
732 Davis, A.; Dean, J.; Devin, M.; Ghemawat, S.; Goodfellow, I.; Harp, A.; Irving, G.; Isard,

- 733 M.; Jia, Y.; Jozefowicz, R.; Kaiser, L.; Kudlur, M.; Levenberg, J.; Mane, D.; Monga, R.;  
734 Moore, S.; Murray, D.; Olah, C.; Schuster, M.; Shlens, J.; Steiner, B.; Sutskever, I.;  
735 Talwar, K.; Tucker, P.; Vanhoucke, V.; Vasudevan, V.; Viegas, F.; Vinyals, O.; Warden,  
736 P.; Wattenberg, M.; Wicke, M.; Yu, Y.; Zheng, X. TensorFlow: Large-Scale Machine  
737 Learning on Heterogeneous Distributed Systems. arXiv March 16, 2016.  
738 <https://doi.org/10.48550/arXiv.1603.04467>.
- 739 (64) Keras 3: A New Multi-Backend Keras, 2023. <https://github.com/keras-team/keras>  
740 (accessed 2023-12-28).
- 741 (65) Asghari, V.; Leung, Y. F.; Hsu, S. C. Deep neural network based framework for complex  
742 correlations in engineering metrics. *Adv Eng Inform.* **2020**, *44*, 101058.
- 743 (66) Wang, R.; Asghari, V.; Hsu, S. C.; Lee, C. J.; Chen, J. H. Detecting Corporate Misconduct  
744 through Random Forest in China's Construction Industry. *J Clean Prod* **2020**, *268*,  
745 122266.
- 746 (67) Corbalá-Robles, L.; Picioreanu, C.; Loosdrecht, M. C. M.; Pérez, J. Analysing the Effects  
747 of the Aeration Pattern and Residual Ammonium Concentration in a Partial Nitrification-  
748 Anammox Process. *Env. Technol* **2016**, *37*, 694–702.
- 749 (68) Skouteris, G.; Arnot, T. C.; Jraou, M.; Feki, F.; Sayadi, S. Modeling Energy Consumption  
750 in Membrane Bioreactors for Wastewater Treatment in North Africa. *Water Env. Res*  
751 **2014**, *86*, 232–244.
- 752 (69) Chen, Z.; Ren, N.; Wang, A.; Zhang, Z. P.; Shi, Y. A Novel Application of TPAD-MBR  
753 System to the Pilot Treatment of Chemical Synthesis-Based Pharmaceutical Wastewater.  
754 *Water Res* **2008**, *42*, 3385–3392.
- 755 (70) Urakami, T.; Sasaki, J.; Suzuki, K. I.; Komagata, K. Characterization and Description of  
756 *Hyphomicrobium Denitrificans* Sp. Nov. *Int J Syst Bacteriol* **1995**, *45*, 528–532.
- 757 (71) Flood, B. E.; Jones, D. S.; Bailey, J. V. *Sedimenticola Thiotaurini* Sp. Nov., a Sulfur-  
758 Oxidizing Bacterium Isolated from Salt Marsh Sediments, and Emended Descriptions of  
759 the Genus *Sedimenticola* and *Sedimenticola Selenatireducens*. *Int J Syst Evol Microbiol*  
760 **2015**, *65*, 2522–2530.
- 761 (72) Collado, L.; Levican, A.; Perez, J.; Figueras, M. J. *Arcobacter Defluvii* Sp. Nov., Isolated  
762 from Sewage Samples. *Int J Syst Evol Microbiol* **2011**, *61*, 2155–2161.
- 763 (73) Chen, Q.; Ni, J.; Ma, T.; Liu, T.; Zheng, M. Bioaugmentation Treatment of Municipal  
764 Wastewater with Heterotrophic-Aerobic Nitrogen Removal Bacteria in a Pilot-Scale  
765 SBR. *Bioresour Technol* **2015**, *183*, 25–32.
- 766 (74) Sorokin, D. Y.; Tourova, T. P.; Lysenko, A. M.; Muyzer, G. Diversity of Culturable  
767 Halophilic Sulfur-Oxidizing Bacteria in Hypersaline Habitats. *Microbiology* **2006**, *152*,  
768 3013–3023.
- 769 (75) Baumann, B.; Snozzi, M.; Zehnder, A. J. B.; Meer, J. R. Dynamics of Denitrification  
770 Activity of *Paracoccus Denitrificans* in Continuous Culture during Aerobic-Anaerobic  
771 Changes. *J Bacteriol* **1996**, *178*, 4367–4374.
- 772 (76) Dou, Q.; Yang, J.; Peng, Y.; Zhang, L. Multipathway of Nitrogen Metabolism Revealed  
773 by Genome-Centered Metatranscriptomics from Pyrite-Guided Mixotrophic Partial  
774 Denitrification/Anammox Installations. *Environ. Sci. Technol.* **2023**, *57* (51), 21791–  
775 21800. <https://doi.org/10.1021/acs.est.3c08192>.

- 776 (77) Narasingarao, P.; Häggblom, M. M. *Sedimenticola Selenatireducens*, Gen. Nov., Sp.  
777 Nov., an Anaerobic Selenate-Respiring Bacterium Isolated from Estuarine Sediment. *Syst*  
778 *Appl Microbiol* **2006**, *29*, 382–388.
- 779 (78) Altamia, M. A.; Shipway, J. R.; Concepcion, G. P.; Haygood, M. G.; Distel, D. L.  
780 *Thiosocius Teredinicola* Gen. Nov., Sp. Nov., a Sulfur-Oxidizing Chemolithoautotrophic  
781 Endosymbiont Cultivated from the Gills of the Giant Shipworm, *Kuphus Polythalamius*.  
782 *Int. J. Syst. Evol. Microbiol.* **2019**, *69* (3), 638–644.  
783 <https://doi.org/10.1099/ijsem.0.003143>.
- 784 (79) Vick, S. H. W.; Tetu, S. G.; Sherwood, N.; Pinetown, K.; Sestak, S.; Vallotton, P.;  
785 Elbourne, L. D. H.; Greenfield, P.; Johnson, E.; Barton, D.; Midgley, D. J.; Paulsen, I. T.  
786 Revealing Colonisation and Biofilm Formation of an Adherent Coal Seam Associated  
787 Microbial Community on a Coal Surface. *Int J Coal Geol* **2016**, *160–161*, 42–50.
- 788 (80) Çinar, Ö.; Hasar, H.; Kinaci, C. Modeling of Submerged Membrane Bioreactor Treating  
789 Cheese Whey Wastewater by Artificial Neural Network. *J Biotechnol* **2006**, *123*, 204–  
790 209.
- 791 (81) Heylen, K.; Vos, P.; Vekeman, B. Draft Genome Sequences of Eight Obligate Methane  
792 Oxidizers Occupying Distinct Niches Based on Their Nitrogen Metabolism. *Genome*  
793 *Announc* **2016**, *4*, 16–17.
- 794 (82) Ogiso, T.; Ueno, C.; Dianou, D.; Huy, T.; Katayama, A.; Kimura, M.; Asakawa, S.  
795 *Methylomonas Koyamae* Sp. Nov., a Type I Methane-Oxidizing Bacterium from  
796 Floodwater of a Rice Paddy Field. *Int J Syst Evol Microbiol* **2012**, *62*, 1832–1837.
- 797 (83) Liesack, W.; Finster, K. Phylogenetic Analysis of Five Strains of Gram-Negative,  
798 Obligately Anaerobic, Sulfur-Reducing Bacteria and Description of *Desulfuromusa* Gen.  
799 Nov., Including *Desulfuromusa Kysingii* Sp. Nov., *Desulfuromusa Bakii* Sp. Nov., and  
800 *Desulfuromusa Succinoxidans* Sp. *Int J Syst Bacteriol* **1994**, *44*, 753–758.
- 801 (84) Xiong, W.; Ye, Y.; He, D.; He, S.; Xiang, Y.; Xiao, J.; Feng, W.; Wu, M.; Yang, Z.;  
802 Wang, D. Dereglulation of Ribosome Biogenesis in Nitrite-Oxidizing Bacteria Leads to  
803 Nitrite Accumulation. *Environ. Sci. Technol.* **2023**, *57* (43), 16673–16684.  
804 <https://doi.org/10.1021/acs.est.3c06002>.
- 805 (85) Jung, Y. T.; Kim, J. H.; Oh, T. K.; Yoon, J. H. *Mariniflexile Aquimaris* Sp. Nov., Isolated  
806 from Seawater, and Emended Description of the Genus *Mariniflexile Nedashkovskaya* et  
807 Al. *Int J Syst Evol Microbiol* **2012**, *62*, 539–544.
- 808 (86) Yang, J.; Chen, Z.; Wang, X.; Zhang, Y.; Li, J.; Zhou, S. Elucidating Nitrogen Removal  
809 Performance and Response Mechanisms of Anammox under Heavy Metal Stress Using  
810 Big Data Analysis and Machine Learning. *Bioresour. Technol.* **2023**, *382*, 129143.  
811 <https://doi.org/10.1016/j.biortech.2023.129143>.

812

813



