



**Environmental
Science**
Water Research & Technology

**Emerging investigator series: Modeling of Wastewater
Treatment Bioprocesses: Current Development and Future
Opportunities**

Journal:	<i>Environmental Science: Water Research & Technology</i>
Manuscript ID	EW-PER-10-2021-000739.R1
Article Type:	Perspective

SCHOLARONE™
Manuscripts

Water Impact

Wastewater treatment bioprocesses are conventionally modeled using mechanistic, data-driven, or hybrid models. Herein, we identify the knowledge gaps in those models. We also propose potential modeling strategies to incorporate genomic data for handling a large amount of the physical, biochemical, and microbiological data collected from biological wastewater treatment systems, with the overarching goal to achieve real-time monitoring and optimize system performance.

1 **Emerging investigator series: Modeling of Wastewater Treatment Bioprocesses: Current**
2 **Development and Future Opportunities**

3

4 Shiyun Yao ^a, Cheng Zhang ^a, Heyang Yuan ^{a,*}

5

6 a. Civil & Environmental Engineering, Department of Civil and Environmental Engineering,
7 Temple University, 1947 N. 12th St, Philadelphia, PA 19122, USA

8

9 Intended for *Environmental Science: Water Research & Technology*

10 Type of contribution: Perspective

11

12 * Corresponding author

13 heyang.yuan@temple.edu

14

15 **Abstract**

16 For more than a half-century, modelers have developed various modeling strategies to facilitate
17 the transition of wastewater treatment bioprocesses from lab-scale demonstration to full-scale
18 applications. This review presents the mathematical fundamentals of mechanistic models,
19 machine learning algorithms of data-driven models, and hybrid modeling strategies for different
20 biological wastewater treatment systems including activated sludge processes, anaerobic
21 digesters, anammox processes, and bioelectrochemical systems. The discussion is focused on the
22 biological principles in those modeling strategies. The conventional Monod expressions are a
23 prevailing tool to describe the mathematical connection between microbial kinetics and state
24 variables in mechanistic models. Stoichiometric equations and steady-state conditions are also
25 required for the mechanistic modeling approach to predict system performance such as the
26 removal of carbon, nitrogen, and phosphorus. On the other hand, data-driven models statistically
27 link the inputs and outputs for the prediction and optimization of system performance with a
28 minimum requirement of a priori knowledge. Although this strategy shows outstanding learning
29 ability of data interpolation, the predictions are often uninterpretable due to the black-box nature.
30 Hybrid modeling strategies have the potential to dress the inherent limitations of standalone
31 models. Currently, the mechanistic and data-driven components in hybrid models are still
32 structured based on microbial kinetics and trained with physical and biochemical data,
33 respectively. This problem can be potentially solved by incorporating genomics data into model
34 construction to link microbial kinetic to microbial population and functional dynamics. We
35 discuss the perspectives of incorporating genomic data into model construction and propose
36 genomics-enabled hybrid modeling strategies for future research.

37

38 **Keywords: hybrid model, wastewater bioprocess, data-driven model, mechanistic model,**
39 **genomics-enabled model**

40 **1. Introduction**

41 A variety of wastewater treatment systems are developed to harness microorganisms to remove
42 organic contaminants and nutrients.¹ These include activated sludge processes, anaerobic
43 digesters, and membrane bioreactors that have been extensively applied in full-scale treatment
44 plants, as well as emerging technologies such as anaerobic ammonium oxidation (anammox) and
45 bioelectrochemical systems. By taking advantage of the metabolic versatility of microorganisms,
46 some of those systems can be engineered to achieve sustainable treatment of complex waste
47 streams. For example, anammox as an anaerobic and autotrophic process converts ammonia to
48 nitrogen gas with little energy input,² and bioelectrochemical systems can recover energy and
49 resource from waste streams.³

50
51 Biological wastewater treatment systems are commonly monitored using real-time biochemical
52 data such as biological oxygen demand (BOD), chemical oxygen demand (COD), and mixed
53 liquor volatile suspended solids (MLVSS) as an indication of biomass concentration.⁴⁻⁶ This is
54 an operationally simple method that leads to a quick assessment of the microbiological activity
55 and system performance. However, those macroscopic state variables sometimes do not help
56 explain the inconsistency in treatment performance caused by environmental and operational
57 perturbations, primarily because they are unable to reflect the complex microbial physiology,
58 community structure, metabolism, and interspecies interactions in the systems.⁴ As conventional
59 bioprocesses face challenges from emerging pollutants and more stringent discharge limits over
60 the past decades,⁷ real-time monitoring and experimental trials are laborious and frequently fall
61 short to provide insights into system optimization.

62

63 The high complexity of bioprocesses and the pressing need to develop more sustainable systems
64 drive environmental researchers and engineers to build computational models to gain
65 mechanistic and predictive understandings of the dynamic behaviors in the systems.⁸ A
66 prevailing strategy is to build mechanistic models, also known as white-box or first-principle
67 models, based on the mathematical expressions of the physical/chemical/biological principles of
68 the processes involved.⁸ This modeling strategy starts with defining the model purpose followed
69 by model structure selection, input identification, data collection/reconciliation, and model
70 calibration/correction.⁹ To improve the robustness of mechanistic models, it is important to
71 understand the fundamental microbial kinetics. Monod expressions are the most well-established
72 tool for modeling microbial growth and substrate utilization kinetics and have been extensively
73 applied to simulate many bioprocesses.¹⁰⁻¹⁵ The kinetic parameters in Monod equations are
74 commonly derived from biochemical measurements.¹⁶ For ill-defined systems that are driven by
75 uncharacterized functional populations, it becomes challenging to determine which biochemical
76 indicators to measure to reliably reflect their growth and substrate utilization kinetics.

77

78 Data-driven models seek to establish statistical connections between the inputs and outputs and
79 require little knowledge about the fundamental principles of the processes involved. This is of
80 great interest to the modeling of bioprocesses, in which the lifestyles of many functional
81 populations remain unknown.¹⁷ Data-driven modeling strategies started to attract
82 environmental/biological engineers' attention in the early 90s with the practical implementation
83 of artificial neural networks (ANNs) and major advances in machine learning.¹⁸⁻²⁰ Using
84 appropriate training datasets and network architecture, neural networks were trained to directly
85 predict wastewater treatment performance,^{20,21} as well as the effects of sludge volume index and

86 total nitrogen concentration on the effluent quality.²² Recent studies demonstrated the
87 applicability of several machine learning algorithms, including support vector machine, random
88 forest, extreme gradient boosting, and k-nearest neighbors, to full-scale digesters.²³ Those
89 algorithms have also been used to train models with biochemical data collected from literature,
90 presenting a powerful tool for data mining.²⁴ Despite the outstanding learning performance, data-
91 driven models are black boxes that often generate uninterpretable predictions.²⁵ This is
92 particularly problematic for bioprocesses when the models are trained solely with physical and
93 biochemical data.^{21-23,26}

94
95 Hybrid models are built through the integration of mechanistic and data-driven sub-models to
96 improve the shortcomings of those individual modeling strategies.²⁷ Hybrid models are ideal for
97 ill-defined systems in which only part of the process can be mechanistically described (e.g., the
98 mass balance of measurable biochemical variables), while another part of the process is too
99 complicated to be derived from first principles (e.g. microbial interactions).²⁸ Depending on the
100 significance of the known and unknown processes in a system, the mechanistic and data-driven
101 sub-models can be integrated through parallel, in series, or a mix of both structures.²⁹ Such
102 integration allows the designer to structure the model more flexibly based on the availability of
103 the a priori knowledge.²⁸ Recent studies demonstrate that hybrid models are more accurate and
104 flexible in predicting the dynamic behaviors of wastewater treatment bioprocesses.^{30,31}

105
106 This review discusses the applications and drawbacks of major modeling approaches that have
107 been developed for simulating wastewater treatment bioprocesses over the past 60 years in
108 chronological order (Figure 1). We first review the mechanistic modeling of several systems

109 including activated sludge processes, anaerobic digesters, anammox, and bioelectrochemical
110 systems. The discussion of this modeling strategy is focused on the biological principles, in
111 particular microbial growth and substrate utilization kinetics. In the following sections, we
112 review data-driven models capable of capturing the dynamic behaviors of bioprocesses, as well
113 as hybrid models that take advantage of both modeling strategies (Figure 2). We also discuss the
114 potential of incorporating high-throughput sequencing data into model construction to unfold the
115 underlying interactions among microbial kinetics, community structure, and functional dynamics.
116 Finally, we propose a conceptual framework to build hybrid models with omics-based results for
117 robust and interpretable prediction of microbial interactions and wastewater treatment
118 performance.

119

120 **2. Mechanistic Modeling of Wastewater Treatment Bioprocesses**

121 **2.1 Early development of microbial kinetic model**

122 Microbial growth is impacted by several factors such as cell metabolic activity, substrate
123 availability, oxygen concentration, temperature, etc. Among them, the substrate is arguably the
124 most important one as it is directly involved in cell metabolism.³² To model the kinetic behaviors
125 of microbial growth, a greater understanding of the effects of substrate concentration on
126 metabolic quotients, specific growth rate, and yield is required. Between the 1920s and 1960s,
127 many hypotheses and models were proposed to explain the biological mechanisms of substrate
128 consumption for cellular metabolism and storage.³³

129

130 The Monod equation is a microbial kinetic model that describes the hyperbolic growth behavior
131 of microbes in batch systems at exponential and steady-state phases.¹⁶ In the Monod equation,

132 the specific growth rate of a culture is a function of the concentration of a given substrate,³⁴ with
133 the substrate saturation constant (i.e., the substrate concentration when the specific growth rate is
134 half of the maximum specific growth rate) indicating the affinity of the substrate to cell growth.¹⁶
135 The Monod equation and its modified expressions were demonstrated to be robust for modeling
136 different pure cultures under varied conditions.³⁵ Since then, a variety of mechanistic models
137 have been built on Monod expressions.^{10,13,36–38}

138
139 Several studies pointed out the inadequacy of the original Monod expression in dealing with
140 substrate inhibition, cell decay, diffusional limitation, etc.^{36,39} Substrate inhibition as a common
141 issue arises when the complex composition in wastewater shows different affinity to microbial
142 cells, leading to competition among functional populations and/or inhibition to cell growth. The
143 Haldane-Andrews equation included an inhibition constant in the Monod expression to reflect
144 substrate inhibition.³⁶ However, because the constant is inferred through the generalized
145 substrate inhibition model, a normal distribution method with an error range, such a modification
146 does not reveal the actual impacts of substrate inhibition on microbial growth.⁴⁰

147

148 **2.2 Modeling of activated sludge processes**

149 The activated sludge models have been developed for more than 70 years since the early 60s.¹¹
150 Early mechanistic models were derived based on steady-state applications in which the cell
151 growth rate remained constant.³⁴ The biochemical parameters used to quantify cell growth were
152 total organic carbon (TOC), total oxygen demand (TOD), and COD along with 5-day BOD.⁴¹
153 They are sufficient for modeling steady-state conditions but have severe limitations in real-time
154 situations where cell growth behaves dynamically under the influence of substrate variation.

155 Oxygen uptake rate and mixed liquor volatile suspended solid were later included to address the
156 limitation.^{36,41} Mathematical techniques such as the feedforward-feedback strategy were also
157 applied to control the flowrate disturbance. These modifications allowed mechanistic models to
158 respond to dynamic situations with real-time biochemical data.⁴¹

159
160 In 1982, the International Association on Water Pollution Research and Control published a
161 preliminary model on activated sludge systems. Later, Dold and Marais incorporated the
162 preliminary model to a final version called the Activated Sludge Model No. 1 (ASM1).^{12,42}
163 ASM1 is considered a reference model and is generally accepted as a fundamental component
164 for wastewater treatment modeling.⁴³ In particular, the Monod equation in ASM1 was proven
165 reasonably appropriate to describe the microbial growth and substrate utilization behaviors in
166 wastewater.⁴⁴

167
168 The microbiological principles of ASM1 include the growth of aerobic and anoxic heterotrophic
169 organisms, the growth of aerobic autotrophic organisms, the decay of heterotrophs and
170 autotrophs, hydrolysis of slowly biodegradable substrate, ammonification, and hydrolysis of
171 organic nitrogen.¹⁰ In addition, the model also describes the following dynamic mass balances
172 that have impacts on biomass concentration: 1) readily biodegradable substrates, 2) slowly
173 biodegradable substrates, 3) inert particulate substances, 4) particulate organic nitrogen, 5)
174 soluble organic nitrogen, 6) ammonia, 7) nitrate, and 8) oxygen.¹⁰ These state variables serve as
175 explicit indicators of the nutrient removal processes. Among all parameters used in ASM1, the
176 growth and decay rates are of key importance as they control the biomass concentration as a
177 function of the influent substrate concentration. The full description of ASM1 and the

178 comprehensive review of model development can be found elsewhere.^{8,10} The robustness of
179 ASM1 has been demonstrated by numerous studies. Forty years after its first implementation,
180 ASM1 is still playing a central role in the mechanistic modeling of bioprocesses and has been
181 incorporated in commercial software as a core structure for the simulation of full-scale
182 wastewater treatment plants.⁴⁵

183

184 Advancing from ASM1, ASM2 incorporates polyphosphate-accumulating organisms (PAO), a
185 functional population enriched during enhanced biological phosphorus removal in activated
186 sludge systems. In addition to common microbial kinetics, the model structure for PAO also
187 includes the storage of glycogen, polyhydroxyalkanoate, and polyphosphate, which are
188 expressed as a function of oxygen availability according to their physiological traits.⁴⁶ In the
189 1990s, ASM2d was developed based on ASM2 by including the ability of PAOs to use internal
190 cellular materials for denitrification,¹⁰ thus linking the metabolism of nitrate and phosphorous
191 under anoxic conditions. In the absence of oxygen, PAOs can use nitrate as a terminal electron
192 acceptor for phosphorous uptake.^{10,47} In ASM2d, a fraction of the maximum growth rate of PAO
193 is assigned to complete denitrification. The fraction varies depending on PAOs' activity
194 including growth, denitrification, and anoxic phosphorus uptake.⁴⁸ ASM2 and ASM2d are
195 comprehensively reviewed by Henze et al.¹⁰

196

197 ASM3 was modified to improve the prediction of oxygen consumption, sludge production,
198 nitrification, and denitrification. Key modifications include cellular storage of organic substrates
199 and the consumption of dead cells through endogenous respiration (instead of the decay and
200 recycling processes described by ASM1).³⁷ With these modifications, ASM3 is more accurate in

201 describing the substrate uptake and storage behaviors, but the enhanced prediction may not be
202 relevant for most treatment plants where ASM1 is sufficient for simulation of general
203 performance.⁴⁶ Therefore, ASM3 is needed only when specific metabolic activities are modeled.

204

205 There are several commercial software packages available to simulate the activated sludge
206 processes following the development of ASM1/2/3. Olsson and Newell provided a detailed
207 overview of the simulator environments for the bioprocesses in wastewater treatment plants.⁴⁹

208 Some of the simulators, such as MATLAB-based Simulink, serve a general purpose with high
209 flexibility to complete simulation. Other simulators contain a library of predetermined models
210 for specific bioprocesses, and the process configuration is a unit-based simulation environment.

211 Examples of this type of simulator are AQUASIM, BioWin, EFOR, GPS-X, SIMBA, STOAT,
212 and WEST. The computation package can solve mechanistic models with multiple mathematical
213 equations simultaneously.⁵⁰ BioWin, for example, is a flexible software tool that includes

214 multiple microbial processes presented in ASMs and anaerobic digester models (ADMs).⁵¹⁻⁵⁵ It
215 has been used to model the biological systems in full-scale wastewater treatment plants,⁵²⁻⁵⁵
216 including activated sludge,⁵⁶ anaerobic digestion,^{57,58} and anammox processes.⁵¹ In some of the

217 studies, BioWin provided a good match between the measured and predicted data (difference
218 <10%) for both small-scale batch reactors and full-scale systems,⁵⁶⁻⁵⁸ and it can perform
219 optimization of sludge retention time and nitrogen removal under different DO and return

220 activated sludge flows.⁵³ In most of the studies, however, model calibration is required to
221 improve the prediction performance. Improper calibration, for example, with nitrogen,
222 phosphorus, and other microbial inhibitory substances has been reported to cause inaccurate

223 prediction of methane production in anaerobic digestion.⁵⁸ Because commercial software still

224 requires both in-depth knowledge of the bioprocess and expertise in modeling, it is not very user-
225 friendly to treatment plant staff.

226

227 **2.3 Modeling of anaerobic digesters**

228 Anaerobic digesters have a highly complex microbial community composed of fermentative
229 bacteria, syntrophs, acidogens/acetogens, and methanogens.⁵⁹ The model structure in early
230 models was improved with a greater understanding of the microbial kinetics of those functional
231 populations and the associated biochemical reactions. For example, Andrews applied the
232 Haldane function on the Monod expression to modify the substrate uptake function with
233 inhibition under rate-limiting conditions,⁶⁰ Andrews and Graef later proposed to include the
234 effects of pH change and buffering through liquid-gas phase interaction and carbonate
235 equilibrium, leading to more accurate modeling of microbial kinetics.⁶¹ Based on those studies,
236 Hill and Barth further included a function that described the inhibition effects of volatile fatty
237 acids (VFAs) and ammonia on methanogens, as well as charge balance to correct temperature-
238 dependent pH.⁶²

239

240 After decades of study of the microbial ecology in anaerobic environments, a task group from the
241 International Water Association consolidated the up-to-date knowledge and formulated
242 Anaerobic Digester Model No. 1 (ADM1) as a common platform model for anaerobic
243 processes.¹³ This model involves 4 typical digestion processes (hydrolysis, acidogenesis,
244 acetogenesis, and methanogenesis) and several physicochemical steps including gas-liquid
245 diffusion, ion association, and dissociation. ADM1 describes in total 29 processes and 32
246 variables at dynamic state, 24 of which are based on the Monod equation and first-order kinetics.

247 The model also allows modifications for specific applications such as sulfate reduction,
248 phosphorous conversion, and mineral precipitation.¹³ ADM1 has been successfully implemented
249 to simulate effluent characteristics with various types of substrates and operating conditions. The
250 most common function of this comprehensive model is to predict and optimize biogas production.
251 Satpathy et al. applied ADM1 to simulate biogas production from rare substrates such as chicken
252 manure.⁶³ Other applications of ADM1 include the prediction of effluent quality from systems
253 fed with winery wastewater⁶⁴ or treating phenol from olive mill waste.⁶⁵ For example, Ozkan-
254 Yucel and Gokcay applied ADM1 to a full-scale anaerobic digester under varying organic
255 loading rates to predict total VFAs and COD in the effluent.⁶⁶ Modified ADM1 can also help
256 troubleshoot operational problems caused by inhibition effects,⁴⁸ in particular, the accumulation
257 of VFAs.⁶⁸

258
259 Despite successful implementation for specific applications, the extensions of sulfate reduction
260 and mineral precipitation can cause a significant computational burden. This is because the
261 precipitation of multiple components (CaHPO_4 , struvite, and other unknown compounds) and
262 their release mechanisms involve a large number of processes, in which the fundamental
263 knowledge is not available for the model to perform *ab initio* prediction.⁶⁹ Another critical issue
264 is that ADM1 is still not able to fully recapitulate the actual functional populations due to the
265 lack of an in-depth understanding of the microbial community. Specifically, ADM1 is structured
266 without considering the production of different short-chain VFAs, alcohols, and hydrogen. The
267 model is thus not able to meet the growing interest in those value-added products.⁶⁴ Shi et al.
268 attempted to solve this problem by redefining the pseudo-stoichiometric dynamic parameters of
269 VFAs and alcohols corresponding to the hydrogen partial pressure.⁷⁰ The modified model

270 successfully predicted the concentrations of acetate, propionate, butyrate, ethanol, and hydrogen
271 with standard errors < 0.04 . However, there remains some discrepancy between the predicted and
272 observed hydrogen levels and effluent COD when the system was fed with high-strength
273 streams.⁷⁰ In terms of methane production, the model did not include the methanogenic
274 population via direct interspecies electron transfer,⁷¹ which was recently found to be a ubiquitous
275 electron transfer mechanism in many engineered and natural environments.⁷² A reaction-
276 diffusion-electrochemistry model composed of activation and ohmic losses predicted that
277 methanogenesis could be an order of magnitude faster via direct interspecies electron transfer
278 than via the classic route of interspecies hydrogen transfer.⁷³ Modified ADM1 further predicted
279 over one-third of the CH_4 produced via this novel electron transfer mechanism,⁷⁴ underpinning
280 its critical role in methane production.

281

282 **2.4 Modeling of emerging bioprocesses**

283 Discovered in the 90s, anammox has been applied as an alternative technology for biological
284 nitrogen removal.⁷⁵ Anammox can be engineered as a single-step process, which is efficient in
285 terms of energy, space, and cost compared with conventional two-step
286 nitrification/denitrification.⁷⁶ In anammox systems, part of the ammonium is oxidized by
287 ammonium oxidizing bacteria through partial nitrification to nitrite, which serves as the electron
288 acceptor for anammox bacteria to oxidize the remaining ammonium to nitrogen gas.⁷⁷ Due to the
289 unique bioprocesses, a previous study built an anammox model from fundamental processes
290 including diffusion, hydrolysis, and microbial kinetics of anammox bacteria for simulation of
291 long-term nitrogen and COD removal by a granular up-flow anaerobic sludge blanket reactor.¹⁵
292 Because the model only assumed cell growth under optimum conditions, the predicted and

293 observed nitrogen removal did not correlate well. The anammox process was later modeled using
294 modified ASM1. Through experimental calibration, the model yielded satisfactory prediction of
295 nitrogen removal efficiency of a laboratory-scale bioreactor.⁷⁸ Interestingly, anammox is better
296 simulated when coupled to other bioprocesses, e.g., sulfur-driven denitrification,¹⁴ as the kinetic
297 parameters in the model can be more accurately estimated by varying both sulfur (sulfite, sulfur,
298 and sulfate) and nitrogen (ammonium, nitrogen, nitrite, and nitrate).¹⁴ Although the models have
299 demonstrated the potential of anammox systems, the modeling of anammox is still challenging
300 because the system performance is highly dependent on the dynamic interactions between
301 ammonium/nitrite oxidizing bacteria and anammox bacteria.¹⁴ A reliable method is to calibrate
302 the maximum growth rates of ammonium oxidizing bacteria, nitrite oxidizing bacteria, and
303 anammox bacteria using experimental data from the full-scale bioreactors and validate the
304 calibration based on the sensitivity analysis.⁷⁹ In addition, the substrate affinity coefficients
305 should be adjusted based on reported literature to fit the microbial substrate utilization under the
306 possible effects of mass transfer in flocs.⁷⁹

307
308 In addition to predicting and optimizing the performance of well-developed bioprocesses such as
309 anammox, mechanistic models can also be implemented for emerging biotechnology such as
310 bioelectrochemical systems to gain an in-depth understanding of their potential and facilitate
311 practical applications. A typical bioelectrochemical system is composed of a cathode and an
312 anode in which the electrochemical reactions are catalyzed by microorganisms. It is challenging
313 to model bioelectrochemical systems primarily because of the close yet uncharacterized
314 connections among microbial kinetics, extracellular electron transfer mechanisms, and
315 electrochemical factors (e.g., internal/external resistance).⁸⁰ A few mathematical models for this

316 novel system have been developed and reviewed elsewhere.^{17,81–83} Noticeably, the growth and
317 substrate utilization rates of the functional groups are still based on the Monod expression and
318 corrected with electron mediator concentration, which in turn is expressed as a function of the
319 substrate utilization rate of electroactive bacterial.^{18,31,84–86} Meanwhile, the Nernst-Michaelis-
320 Menten equation was used to calculate the electron transfer rate in the system⁸⁷ and the Nernst-
321 Planck equation was introduced to represent ion diffusion through the membrane between anode
322 and cathode.^{88,89}

323

324 **2.5 Inherent drawbacks of mechanistic modeling**

325 The model structure in response to intracellular biochemical reactions is oftentimes inadequate,
326 which represents a major limitation of mechanistic models. Structural inadequacy stems from the
327 generalization of microbial growth and substrate utilization kinetics with Monod expressions.⁸
328 Ideally, the Monod expression for the growth rate of a functional population should consider all
329 degradable and inhibitory compounds and the corresponding factors. The substrates in practical
330 applications (e.g., wastewater) contain various degradable and inhibitory compounds that result
331 in highly complex microbial communities, in which many of the functional populations are
332 uncharacterized.^{90–92} Therefore, it is not clear which substrates should be experimentally
333 measured to reflect their growth and substrate utilization kinetics. The absence of expressions for
334 novel methanogenesis mechanism via direct interspecies electron transfer is an example of
335 inadequate model structure.⁵⁷ Such an issue is more problematic for emerging bioprocesses. The
336 models for bioelectrochemical systems were typically structured with fermentative, electroactive,
337 and methanogenic populations. Although the microbial community in bioelectrochemical
338 systems is much more diverse than those in anaerobic digestors,^{93,94} uncharacterized populations

339 cannot be incorporated into the model structure because their physiological traits related to
340 growth and substrate utilization are poorly understood.^{95–97}

341
342 Another challenge lies in the unmeasurable kinetic parameters in Monod expressions. Most of
343 the kinetic parameters such as inhibition constant, maximum growth rates, half-saturation
344 constants, and substrate utilization rates can only be derived from biochemical measurements
345 (e.g., substrate concentration).⁹⁸ From an experimental perspective, when microbial kinetics
346 responds to changing influent characteristics in real-time, it is unrealistic to derive the kinetic
347 parameters throughout the variation. As the operating condition changes, biochemical
348 measurements can vary significantly, and the estimation of kinetic parameters becomes
349 conditional.²⁷ The identifiability of unmeasurable kinetic parameters, i.e., the possibility to
350 derive a unique set of values for the parameters from experimental data, is then of great concern.
351 Firstly, uncertainty issues arise because the derivation of the kinetic parameters is obtained
352 through rate-controlling experiments with specific temperatures, pH, BOD, and COD, etc.,³⁴
353 while in wastewater treatment processes, microbial cells are exposed to varied substrates. Those
354 parameters need constant calibration. An uncalibrated model is not likely to yield accurate
355 predictions. As a result, it is not certain to what extent the predictions can be used to explain the
356 observed physical, chemical, and biological mechanisms..⁹⁹ Secondly, to identify the kinetics
357 parameters, the identification of the biochemical and physical parameters responsible for
358 biomass concentrations must be achieved first. Flotats et al. applied the Taylor Series Expansion
359 to four state variables (acetate, propionate, valerate, and methane) of ADM1 to identify the
360 parameters related to valerate consumption and biomass concentration.¹⁰⁰ Such a method has
361 only been reported to solve simple models with a few unidentifiable parameters.^{101,102} For

362 complex models such as second-order models, parameter identification was handled with the
363 asymptotic behavior of the maximum likelihood estimator and multiple shooting approach
364 described in Muller et al.¹⁰³ Those previous studies collectively show that the identification and
365 computation processes of kinetic parameters are cumbersome and uncertain in reality.

366

367 **3. Data-Driven Modeling of Wastewater Treatment Bioprocesses**

368 **3.1 Neural networks**

369 Neural networks, first reported in 1943, are arguably the most prevailing data-driven models
370 across various research fields.^{104,105} A neural network is composed of multiple layers of
371 interconnected nodes (neurons), through which the inputs are propagated to the final output
372 layer.¹⁰⁶ Each input to a neuron has a weight factor that determines the interconnection strength
373 to the next neuron. By adjusting the weight factors, a neural network can be properly trained to
374 perform problem-solving. The training algorithm can be divided into three types: supervised,
375 unsupervised, and hybrid training. In supervised training, neural networks are trained with a
376 labeled dataset that provides feedback about the prediction accuracy.¹⁰⁶ Unsupervised training
377 allows networks to be trained with unlabeled data, and the algorithm extracts features and
378 patterns on its own.¹⁰⁶ The hybrid training strategy uses unsupervised training for the hidden
379 neurons and supervised training for the output neurons.¹⁰⁷

380

381 Neural networks were first implemented for continuously stirred bioreactors to predict the
382 fermentation products and pH in the effluent.¹⁰⁸ Boger applied the modeling strategy to full-scale
383 wastewater treatment plants and showed that neural networks could be a solution for simulating
384 expert rules, a set of boundary values that confined the neural network prediction, from historical

385 operating data.¹⁹ Traditionally, neural networks were trained in a feed-forward fashion, meaning
386 that the feed was directed forward-only through layers of training. Yang and Linkenst found that
387 the back-propagation method, which fed outputs from a random layer back to a previous layer,
388 could help the network lower the error rate of prediction.¹⁰⁹ Backpropagation thus became a
389 prevailing algorithm of neural network modeling. Several studies used this method to predict
390 effluent COD, biogas production, and NH_4^+ -N removal in different biological wastewater
391 treatment systems including activated sludge processes,²¹ up-flow anaerobic sludge blanket
392 reactors,²⁶ sequencing batch reactors,¹¹⁰ and anaerobic digesters.¹¹¹

393
394 Previous studies have combined neural networks with other types of data-driven models to
395 improve the simulation of ill-defined systems. One of the strategies is to use the genetic
396 algorithm to select the initial dataset for downstream neural network training, thereby identifying
397 the optimal training parameters and reducing the computational burden.¹¹² Bagheri et al. have
398 successfully applied the genetic algorithm to optimize the weights and thresholds of neural
399 networks to accurately predict sludge volume index.²² Neural networks can also be coupled to
400 particle swarm optimization, a population-based optimization technique that searches for optimal
401 weights and biases through multiple iterations of particle positions in a given search space. This
402 coupling approach was intended to lower the training time and computational cost for finding an
403 optimal neural network structure.¹¹³ Compared to the genetic algorithm, neural networks coupled
404 to particle swarm optimization is more memory-efficient in searching optimal weight parameters
405 but is less practical because it has no crossover and mutation in its operator.¹¹⁴

406

407 Neural networks have been proven feasible for interpolation within the training data range.⁹¹
408 Insufficient training data can result in low prediction accuracy, particularly when prediction is
409 performed with a dataset of fewer than 10 samples.⁹³ Another issue is that the model structure is
410 often determined based on a trial-and-error approach, leading to significant time consumption
411 and computational cost.⁹² The fact that neural networks are black box built based on data fitting
412 rather than the mechanistic understanding of the processes suggests that their outputs cannot be
413 used to explain the mechanisms where the inputs are sourced.²⁵

414

415 **3.2 Random forest**

416 Random forest is initially developed as a stochastic discrimination approach for classification
417 purposes in the 90s.¹¹⁵ Later, the approach was extended to combine bagging and random
418 selection features to construct a collection of decision-making trees with control variance.¹¹⁶ To
419 use the random forest algorithm, the input data are classified through layers of tree branches
420 consisting of a variety of features and classes, and multiple trees composed of the same number
421 of features and classes are collectively used for prediction.¹¹⁶

422

423 In the wastewater treatment field, the random forest algorithm has been implemented for
424 activated sludge processes, anaerobic digesters, membrane bioreactors, and anammox
425 processes.^{117–120} The main use of random forest-based models includes the prediction of system
426 performance, fault finding, big data handling, model comparisons, and exploration of datasets
427 with applicable reservations and constraints.¹²¹ Although random forest-based models, similar to
428 other data-driven models, are not able to integrate biological principles, these models allow for
429 the identification of key features and conditions that are most influential on the process. The

430 inference can shed light on the underlying biochemical mechanisms. Song et al. implemented
431 this modeling approach with wastewater treatment inputs as multivariate datasets to predict N₂O
432 emission from the aerated zones of activated sludge processes.¹¹⁸ Based on the model inference,
433 they identified inorganic carbon concentration and specific ammonia oxidation activity as two of
434 the dominant factors that determined treatment performance.¹¹⁸ The model was further used to
435 identify the different mechanisms of N₂O generation in oxic and anoxic environments and
436 demonstrated the key role of N₂O in those zones in promoting niche-specific biochemical
437 reactions.¹²²

438

439 Random forest can be combined with other algorithms such as principal component analysis
440 (PCA) and neural networks to improve the prediction of effluent quality. Preprocessing of data
441 using PCA could enhance the robustness of random forest-based models, leading to a more
442 accurate prediction of membrane flux in membrane bioreactors as compared with neural
443 networks.¹¹⁷ When coupled to neural networks, random forest-based models could be trained to
444 predict the settleability in the biological reactor chamber,¹¹⁹ as well as to evaluate the effects of
445 key operating factors on treatment performance.¹²³ The results suggested that such a combined
446 strategy could help achieve real-time monitoring and optimize operating conditions.¹²⁴

447

448 **3.3 Fuzzy logic**

449 Fuzzy logic is an if-then algorithm that can be used to develop a set of flexible rules for
450 diagnosis and control. A fuzzy logic-based system has four robust components: a fuzzifier, a
451 fuzzy rule-base, an inference engine, and a defuzzifier.¹²⁵ The fuzzifier is responsible for
452 converting crisp inputs into fuzzy sets, which are mapped by the inference engine to produce

453 another fuzzy set as the output. The fuzzy rule base is a collection of rules that guide the fuzzy
454 engine to produce the outputs. Finally, the defuzzifier transfers the output fuzzy sets back to
455 crisp values.¹²⁵

456

457 Fuzzy logic can be applied to numerous scenarios in wastewater treatment bioprocesses such as
458 diagnosis and control of sequencing batch reactor processes,¹²⁶ simulations and prediction of
459 phosphorus removal,¹²⁷ as well as design, evaluation, and decision optimization of activated
460 sludge processes.¹²⁸ Robles et al. developed a fuzzy logic-based controller to optimize biogas
461 production and VFA concentration at varied influent flow rates.¹²⁹ The designed controller was
462 able to help prevent acidification in a closed-loop setting.¹²⁹

463

464 A combination of fuzzy logic and neural networks could bring together the learning powers of
465 both algorithms, enabling fault tolerance during the modeling of complex systems.^{107,108} For
466 example, the adaptive neuro-fuzzy inference system that pairs neural networks with fuzzy logic
467 allows modelers to insert *a priori* knowledge into the neural network structure as rules for
468 training. The combined modeling strategy was used to predict suspended solids, COD, pH and
469 DO levels in activated sludge systems.¹⁰⁸⁻¹¹⁰ Using the adaptive neuro-fuzzy inference system,
470 Essienubong et al. obtained a strong correlation between the experimental and predicted biogas
471 production with temperature, pH, substrate/water ratio, and hydraulic retention time as the
472 inputs.¹³⁴ The work by Hosseinzadeh et al. also demonstrated a higher sensitivity of the adaptive
473 neuro-fuzzy inference system when predicting water flux in an osmotic membrane bioreactor.¹¹²
474 In addition to neural networks, fuzzy logic was coupled to genetic algorithms and particle swarm

475 optimization, which outperformed the adaptive neuro-fuzzy inference system during the
476 prediction of BOD, ammonium, and suspended solids in specific bioprocesses.¹³⁶

477
478 Although fuzzy logic models are powerful tools for predicting system outputs using observable
479 environmental data and human-like logic, these data-driven models cannot capture the behaviors
480 of the complex kinetic reactions in engineered biological bioprocesses, which presents the major
481 criticism among other concerns of implementation.¹³⁷ From an engineering point of view,
482 detailed descriptions of the chemical, physical, and microbiological principles in bioprocesses
483 and computing-based predictive methodology are equally important.¹³⁸ Unfortunately, fuzzy
484 logic systems, like most of the data-driven models, are incapable of providing mechanistic
485 insight into troubleshooting and system optimization due to their black-box nature.

486

487 **4. Hybrid Models That Address the Limitations of Conventional Modeling Strategies**

488 The concept of combining mechanistic and data-driven sub-models for hybrid modeling was first
489 proposed in the early 90s and immediately implemented for a fermentation bioprocess to reduce
490 the dependence on microbial kinetics.²⁹ Hybrid modeling strategies were further examined with
491 activated sludge processes and anaerobic digesters through parallel or serial combinations of the
492 mechanistic and data-driven components.^{139,140}

493

494 In a parallel structure, the outputs from the mechanistic and data-driven components are
495 combined primarily through pure superposition (i.e., summation of the outputs).^{106,140–142}
496 Weighing functions can be introduced to adjust the weight of the outputs, thereby improving the
497 overall prediction accuracy.¹⁴³ It should be noted that the prediction performance of a parallel-

498 structured hybrid model is highly dependent on the robustness of the individual sub-models.¹⁴¹ In
499 cases when a biological system is too dynamics/nonlinear and some of the biochemical data are
500 too expensive to collect in real-time (e.g., H₂ in anaerobic digesters), neither the mechanistic nor
501 the data-driven components could predict accurately, leading to poor prediction performance of
502 the hybrid model. It is therefore argued that the sub-models in a parallel structure are not well
503 integrated due to the lack of interactions (e.g., cross-feeding of the outputs as done in a serial
504 structure).

505

506 In a serial structure, the data-driven component acts as a parameter simulator and estimates
507 kinetic parameters for the mechanistic component to complete simulation.^{29,140} Serial coupling of
508 the sub-models leverages the prediction power of the data-driven component and is sometimes
509 capable of extrapolating system outputs outside of the observation range. A mixture of both
510 structures has been implemented for chemical systems but not bioprocesses.¹⁰⁶ As the selection
511 of the structures depends on the availability of the mechanistic information,¹⁴² we argue that with
512 more functional populations being characterized in wastewater biological systems,¹⁴⁴ a serial or
513 mixed structure may better reflect the underlying biological mechanisms whilst accurately
514 predicting the system performance.

515

516 One of the advantages of hybrid modeling is that unmeasurable mechanistic parameters, in
517 particular microbial growth and substrate utilization rates, can be determined by the data-driven
518 component (Figure 2).^{29,145} This was done in the first study of hybrid modeling of engineered
519 bioprocesses, in which a neural network was trained to estimate the specific growth rate of the
520 overall microbial community, and the outputs were used to establish a biomass balance, resulting

521 in more accurate predictions than those from a standalone neural network.²⁹ Compared to
522 conventional mechanistic models that perform one-time estimation of the kinetic parameters, the
523 data-driven component in hybrid models allows the microbial kinetics to be updated in a timely
524 manner based on the collected data, thereby making the hybrid model more robust under varied
525 conditions.²⁹ Another benefit is that the data-driven component can capture more dynamic data
526 to compensate for the prediction of the mechanistic component. For instance, a neural network
527 was used to estimate the operational data at time $t+1$ based on the influent and operational data at
528 time t , and the outputs were subsequently used to correct the mechanistic predictions at time
529 $t+1$.¹⁴⁰ In this way, the hybrid model could capture the disturbance caused by shock loadings of
530 toxic compounds and deliver more accurate prediction of the effluent composition.¹⁴⁰

531
532 Similar to standalone mechanistic models, hybrid models still contain a Monod expression-based
533 model structure. As previously discussed, Monod expressions approximate the physiological
534 traits of functional populations (i.e., microbial growth and substrate utilization) by assuming a
535 homogeneous culture.^{16,35} Real wastewater and sewage sludge are highly heterogeneous with
536 various degradable and inhibitory compounds that result in diverse microbial populations.^{43,55}
537 For well-characterized populations with known functions, unmeasurable kinetic parameters are
538 predominantly derived from biochemical measurements.^{14,16,34,35,40} Although the data-driven
539 component can help correct the estimates and improve the prediction performance, the estimated
540 kinetic parameters do not necessarily reflect the actual activity of those populations. This is even
541 more problematic for uncharacterized functional populations, whose microbial activity
542 information can be inferred with the data-driven component but does not help interpret the final
543 prediction because those populations are frequently overlooked in the mechanistic component.

544 Therefore, conventional hybrid models constructed with biochemical and physical data still
545 suffer from interpretation issues. A potential solution is to incorporate microbial population and
546 functional dynamics directly into the mechanistic and data-driven components.

547

548 **5. Genomics-Enabled Modeling Strategies for Accurate and Interpretable Prediction**

549 **5.1 Genomics-enabled data-driven modeling**

550 The rapid development of high throughput sequencing techniques and bioinformatics has led to a
551 greater understanding of the microbiomes in wastewater treatment bioprocesses.^{146–149}
552 Recovering the 16s rRNA sequences allows us to unfold the taxonomy and phylogeny of the
553 core populations in activated sludge, anaerobic digester, and many other systems.¹⁴⁶ Meanwhile,
554 metagenomic and meta-transcriptomic data have greatly advanced our knowledge about the
555 genetic potential and functional dynamics of uncharacterized populations.^{147,148} The findings
556 gained with those powerful tools have validated the mechanistic structure of existing models
557 formulated with known functional populations.¹⁴⁹ The million-dollar question now is how to
558 incorporate those genomic data into modeling in a more direct manner (Figure 3).

559

560 Two pioneering studies integrated 16S rRNA amplicon sequencing data with machine learning
561 algorithms (neural networks and Bayesian networks) to reconstruct the microbial communities in
562 natural ecosystems.^{150,151} Following similar strategies, several machine learning-based models
563 were trained with microbial taxon abundance and generated semi-interpretable predictions of
564 system performance and stability.^{152,153} The inferences suggested that for specific systems,
565 classifying taxonomic data at the family level could enhance the prediction accuracy, whilst
566 abundances of specific genera could act as better predictors, highlighting the potential to improve

567 the prediction interpretability with proper data preparation. Before training Bayesian networks,
568 Yuan et al. prepared the genomic data collected from a bioelectrochemical system by selecting
569 dominant taxa at the phylum, genus (Figure 4A), and operational taxonomic unit levels.¹⁵⁴ The
570 genomics-enabled data-driven modeling approach was rigorously cross-validated using three
571 validation strategies.¹⁵⁴ Firstly, the difference between the predicted and observed relative
572 abundances of the selected populations remained within an acceptable range as indicated by a
573 relative root-mean-square error (RMSE) of ~20%. Secondly, the microbial communities
574 reconstructed with the predicted abundances of the selected populations shared high Bray-Curtis
575 similarity with the observed communities at all taxonomic levels. Finally, the predicted system
576 outputs agreed well with the experimental data. For example, current production as the most
577 important system performance for bioelectrochemical systems was predicted with high accuracy
578 at the order level ($R^2 = 0.97$ for prediction vs. observation, Figure 4B). After validation, the
579 model was used to predict current production as a function of operating conditions (e.g.,
580 substrate salinity, Figure 4C) and provided practical insights into system optimization.

581
582 Functional genomic data as the training input can improve prediction interpretability. A previous
583 study trained ANNs with gene expression levels to infer metabolic behaviors, resulting in a
584 plausible explanation of microbes' stress adaptation behaviors under environmental
585 perturbations.¹⁵⁵ Using a similar but more dynamic modeling strategy, Yuan et al. trained
586 Bayesian networks with meta-transcriptomic data to explain the contribution of interspecies
587 hydrogen transfer and direct interspecies electron transfer to methanogenesis.³⁰ It is highly
588 desired to develop a predictive understanding of the involvement of the two mechanisms due to
589 the lack of measurement techniques.^{73,74} To prepare data for model training, the genes for alcohol

590 metabolism, hydrogen metabolism, extracellular electron transfer, and methanogenesis were
591 extracted from the metagenome-assembled genomes of the dominant microbes. A Bayesian
592 network trained with those genes is composed of two components (Figure 5A): upstream gene-
593 gene interactions that predict the expression level of the relevant genes in methanogens, and a
594 downstream sub-network that links the genes encoding methanogenesis to methane. A complete
595 network could accurately predict methane production ($R^2 = 0.96$ for prediction vs. observation,
596 Figure 5B). To statistically infer the contribution of the electron transfer mechanisms, relevant
597 genes were manually silenced. When the simulation was performed without the genes for
598 hydrogen metabolism, the prediction accuracy was significantly compromised as evidenced by
599 the noticeable difference between the predicted and observed methane production and high
600 RMSE (Δ IHT in Figure 5B). In contrast, the prediction remained accurate with *in silico* knockout
601 of the genes for direct interspecies electron transfer (Δ DIET in Figure 5B). The results thus
602 implied a more critical role of hydrogen-mediated electron transfer in methane production.

603

604 **5.2 Genomics-enabled hybrid modeling**

605 Thus far, genomic data have only been used to train data-driven models for semi-interpretable
606 prediction. There is a growing interest to incorporate it into hybrid modeling to predict the
607 underlying mechanisms for system design and optimization. One potential strategy is to infer
608 unidentifiable kinetic parameters from microbial population and gene dynamics. The mechanistic
609 component can be formulated following conventional modeling procedures to estimate kinetic
610 parameters, and the estimates can be combined with microbial taxon abundance and operating
611 conditions to train the data-driven component. The resulting hybrid model is therefore a kinetics

612 simulator that statistically infers kinetic parameters, which can then be fed back to the
613 mechanistic component for prediction of system performance.

614

615 The concept has been proven valid by a recent study,³¹ in which a genomics-enabled hybrid
616 model was implemented for bioelectrochemical systems based on 77 samples collected from 13
617 publications. The mechanistic component of the hybrid model was built to estimate the
618 maximum growth and substrate utilization rates of three functional populations: fermentative,
619 electroactive, and methanogenic microbes, which were subsequently combined with the relative
620 abundances of 38 core taxa at the genus level to train a hybrid Bayesian network (Figure 6A).
621 When examined with six new samples that were not included in network training, the hybrid
622 model achieved the most accurate prediction of current production (Hybrid + Mechanistic in
623 Figure 6B) compared with standalone data-driven models. The enhanced prediction performance
624 of the hybrid model likely results from the close connection between population dynamics and
625 microbial kinetics.

626

627 An alternative strategy to incorporate genomic data into hybrid modeling is to replace Monod
628 expression-based model structures in conventional mechanistic models by simulating microbial
629 population dynamics. This can be achieved with an iterative strategy, in which the data-driven
630 component trained with processed genomic data infers instantaneous biochemical and microbial
631 intermediates, and the intermediates are fed into the mechanistic component to predict steady-
632 state biochemical outputs. This novel strategy mimics microbial community assembly driven by
633 environmental perturbations in engineered systems:^{156–158} operating conditions and biochemical
634 inputs impose selective pressure and together shape the microbial community structure, and the

635 enriched functional populations produce biochemical intermediates that are rapidly mixed with
636 the inputs to form a new steady-state, causing the community structure to further shift until the
637 biochemical outputs reach equilibrium. The proposed strategy thus frees models from microbial
638 kinetics-based structures, while the abiotic and microbial processes and their interplay are
639 revealed in each iteration. Successful modeling with this strategy relies on the ability of hybrid
640 models to act as a community simulator to predict microbial taxon abundance and reconstruct
641 microbial communities *in silico*. For example, Bayesian networks could infer the relative
642 abundances of dominant taxa, resulting in Bray-Curtis similarity of over 90% between the
643 simulated and observed microbial communities at the phylum level.¹⁵⁶ However, the similarity
644 dropped to 83% at the order level and 69% at the OTU level, likely because of the presence of
645 functionally redundant taxa in the small data pool. The potential of the proposed strategy
646 warrants investigation with big data collected from global databases.

647

648 **6. Conclusions**

649 This review focuses on three major types of models: mechanistic, data-driven, and hybrid models.
650 Mechanistic models can provide fundamental insights but need laborious calibration because the
651 Monod-based model structure is inadequate to reflect the biological principles whilst the
652 microbial kinetic parameters are largely unidentifiable. As a result, a mechanistic model built for
653 a specific system frequently falls short when applied to other bioprocesses. Data-driven models
654 can provide predictive insights but yield uninterpretable predictions due to their black-box nature.
655 Hybrid models are believed to overcome the issues of structural inadequacy, parametric
656 unidentifiability, and uninterpretable prediction of the standalone models. Recent
657 biotechnological development such as high throughput sequencing data and omics-based analysis

658 can further enable the incorporation of microbial population and functional dynamics into the
659 model to directly reflect the biological principles. Genomics-enabled hybrid modeling strategies
660 require the mechanistic and data-driven components to be integrated interactively. Two strategies
661 are proposed: kinetics simulator and community simulator, and their applicability warrant further
662 studies. Although hybrid models can potentially overcome the drawbacks of standalone models,
663 the main rationale of modeling selection and design is largely dependent on its intended use.
664 Additionally, the availability of omics-based data and computational cost require more effort in
665 preparation, collection, process, and analysis, which demands technical labor, time, and financial
666 investment. All these factors need to be taken into consideration when modelers design and
667 modify existing models.

668

669 **Conflicts of Interest**

670 There are no conflicts of interest to declare.

671

672 **Acknowledgment**

673 This work was supported by the U.S. Department of Agriculture [Award No. 2020-67019-
674 31027].

675 **Reference**

- 676 1 USEPA, *Emerging Technologies for Wastewater Treatment and In-Plant Wet Weather*
677 *Management*, U.S. Environmental Protection Agency, Fairfax, 2013.
- 678 3 J. S. Seelam, S. A. Maesara, G. Mohanakrishna, S. A. Patil, A. Ter Heijne and D. Pant, in
679 *Waste Biorefinery: Potential and Perspectives*, Elsevier, 2018, pp. 535–570.
- 680 4 *Module 16: the Activated sludge process part II*, Pennsylvania Department of Environmental
681 Protection, Harrisburg, 2014.
- 682 5 *Biological Treatment: suspended growth processes study guide*, Wisconsin Department of
683 Natural Resources, Madison, WI, 2015.
- 684 6 R. Snyder and D. Wyant, *Activated sludge process control: training manual for wastewater*
685 *treatment plant operators*, State of Michigan Department of Environmental Quality.
- 686 7 A. Xu, Y.-H. Wu, Z. Chen, G. Wu, Q. Wu, F. Ling, W. E. Huang and H.-Y. Hu, Towards the
687 new era of wastewater treatment of China: Development history, current status, and future
688 directions, *Water Cycle*, 2020, **1**, 80–87.
- 689 8 U. Jeppsson, *Modelling Aspects of Wastewater Treatment Processes*, 1996.
- 690 7 K. V. Gernaey, M. C. M. Van Loosdrecht, M. Henze, M. Lind and S. B. Jørgensen,
691 Activated sludge wastewater treatment plant modelling and simulation: State of the art, 2004,
692 **19**, 763-783.
- 693 10 M. Henze, W. Gujer, T. Mino, T. Matsuo, M. Wentzel, G. Marais and M. Loosdrecht,
694 Activated sludge model No.2D, ASM2D, *Water Sci. Technol.*, 1999, **39**, 165–182.
- 695 11 M. C. M. Van Loosdrecht, C. M. Lopez-Vazquez, S. C. F. Meijer, C. M. Hooijmans and D.
696 Brdjanovic, Twenty-five years of ASM1: Past, present and future of wastewater treatment
697 modelling, *J. Hydroinformatics*, 2015, **17**, 697–718.
- 698 12 M. Henze, W. Gujer, T. Mino and M. Van Loosdrecht, Activated sludge models ASM1,
699 ASM2, ASM2d and ASM3, *Water Intell. Online*, 2006, **5**, 47–77.
- 700 13 D. J. Batstone, J. Keller, I. Angelidaki, S. V Kalyuzhnyi, S. G. Pavlostathis, A. Rozzi, W. T.
701 M. Sanders, H. Siegrist and V. A. Vavilin, *The IWA Anaerobic Digestion Model No 1*
702 *(ADMI)*, 2002.
- 703 14 Y.-F. Deng, W.-T. Tang, H. Huang, J. Qian, D. Wu and G.-H. Chen, Development of a
704 kinetic model to evaluate thiosulfate-driven denitrification and anammox (TDDA) process,
705 *Water Res.*, 2021, **198**, 117155.
- 706 15 B. J. Ni, Y. P. Chen, S. Y. Liu, F. Fang, W. M. Xie and H. Q. Yu, Modeling a granule-based
707 anaerobic ammonium oxidizing (ANAMMOX) process, *Biotechnol. Bioeng.*, 2009, **103**,
708 490–499.
- 709 16 J. Monod, The Growth of Bacterial Cultures, *Annu. Rev. Microbiol.*, 1949, **3**, 371–394.
- 710 17 S. Luo, H. Sun, Q. Ping, R. Jin and Z. He, A Review of Modeling Bioelectrochemical
711 Systems: Engineering and Statistical Aspects, *Energies*, 2016, **9**, 1–27.
- 712 18 Y. Liu, M. Qin, S. Luo, Z. He and R. Qiao, Understanding Ammonium Transport in
713 Bioelectrochemical Systems towards its Recovery, *Sci. Rep.*, 2016, **6**, 1–10.
- 714 19 Z. Boger, Application of neural networks to water and wastewater treatment plant operation,
715 *ISA Trans.*, 1992, **31**, 25–33.
- 716 20 N. Bhat and T. J. McAvoy, Use of neural nets for dynamic modeling and control of chemical
717 process systems, *Comput. Chem. Eng.*, 1990, **14**, 573–582.

- 718 21 H. Moral, A. Aksoy and C. F. Gokcay, Modeling of the activated sludge process by using
719 artificial neural networks with automated architecture screening, *Comput. Chem. Eng.*, 2008,
720 **32**, 2471–2478.
- 721 22 M. Bagheri, S. A. Mirbagheri, Z. Bagheri and A. M. Kamarkhani, Modeling and
722 optimization of activated sludge bulking for a real wastewater treatment plant using hybrid
723 artificial neural networks-genetic algorithm approach, *Process Saf. Environ. Prot.*, 2015, **95**,
724 12–25.
- 725 23 D. De Clercq, D. Jalota, R. Shang, K. Ni, Z. Zhang, A. Khan, Z. Wen, L. Caicedo and K.
726 Yuan, Machine learning powered software for accurate prediction of biogas production: A
727 case study on industrial-scale Chinese production data, *J. Clean. Prod.*, 2019, **218**, 390–399.
- 728 24 L. Wang, F. Long, W. Liao and H. Liu, Prediction of anaerobic digestion performance and
729 identification of critical operational parameters using machine learning algorithms,
730 *Bioresour. Technol.*, 2020, **298**, 1–9.
- 731 25 C. Rudin, Stop explaining black box machine learning models for high stakes decisions and
732 use interpretable models instead, *Nat. Mach. Intell.*, 2019, **1**, 206–215.
- 733 26 C. Mendes, R. da Silva Magalhes, K. Esquerre and L. M. Queiroz, Artificial neural network
734 modeling for predicting organic matter in a full-scale up-flow anaerobic sludge blanket
735 (UASB) reactor, *Environ. Model. Assess.*, 2015, **20**, 625–635.
- 736 27 M. J. Wade, Not just numbers: Mathematical modelling and its contribution to anaerobic
737 digestion processes, *Processes*, 2020, **8**, 1-31, DOI:10.3390/PR8080888.
- 738 28 T. Bohlin, *Practical grey-box process identification: theory and applications*, Springer,
739 London, 2006.
- 740 29 D. C. Psichogios and L. H. Ungar, A hybrid neural network-first principles approach to
741 process modeling, *AIChE J.*, 1992, **38**, 1499–1511.
- 742 30 H. Yuan, X. Wang, T.-Y. Lin, J. Kim and W.-T. Liu, Disentangling Syntrophic Electron
743 Transfer Mechanisms in Methanogenesis Through Electrochemical Stimulation, Omics, and
744 Machine Learning, 2021, **11**, 1-22, DOI:10.21203/rs.3.rs-288821/v1.
- 745 31 Z. Cheng, S. Yao and H. Yuan, Linking Population Dynamics to Microbial Kinetics for
746 Hybrid Modeling of Engineered Bioprocesses, *bioRxiv*, 2021, 2021.04.15.440059.
- 747 32 B. E. Rittmann and P. L. McCarty, *Environmental biotechnology: Principles and*
748 *applications*, McGraw-Hill Education, 2001.
- 749 33 J. E. McKee, G. M. Fair and L. S. Kraus, Load Distribution in the Activated Sludge Process,
750 *Sew. Works J.*, 1942, **14**, 121–146.
- 751 34 K. Kovárová-Kovar and T. Egli, Growth kinetics of suspended microbial cells: from single-
752 substrate-controlled growth to mixed-substrate kinetics, *Microbiol. Mol. Biol. Rev.*, 1998, **62**,
753 646–666.
- 754 35 J. Heijnen and B. Romein, Derivation of Kinetic Equations for Growth on Single Substrates
755 Based on General Properties of a Simple Metabolic Network, 1995, **11**, 712–716.
- 756 36 J. F. Andrews, A mathematical model for the continuous culture of microorganisms utilizing
757 inhibitory substrates, *Biotechnol. Bioeng.*, 1968, **10**, 707–723.
- 758 37 W. Gujer, M. Henze, T. Mino and M. Van Loosdrecht, Activated sludge model no. 3, *Water*
759 *Sci. Technol.*, 1999, **39**, 183–193.
- 760 38 J. Alex, L. Benedetti, J. Copp, K. V. Gernaey, U. Jeppsson, I. Nopens, M. N. Pons, L. Rieger,
761 C. Rosen, J. P. Steyer, P. Vanrolleghem and S. Winkler, Benchmark Simulation Model no. 1
762 (BSM1), 2008, 62.

- 763 39 M. Mohammadi, A. R. Mohamed, G. D. Najafpour, H. Younesi and M. H. Uzir, Kinetic
764 Studies on Fermentative Production of Biofuel from Synthesis Gas Using *Clostridium*
765 *ljungdahlii*, *Sci. World J.*, 2014, **2014**, 1–9.
- 766 40 G. C. Okpokwasili and C. O. Nweke, Microbial growth and substrate utilization kinetics, *Afr.*
767 *Journal Biotechnol.*, 2005, **5**, 305–317.
- 768 41 M. K. Stenstrom, W. Kido, R. F. Shanks and M. Mulkerin, Estimating Oxygen Transfer
769 Capacity of a Full-Scale Pure Oxygen Activated Sludge Plant, *J. Water Pollut. Control Fed.*,
770 1989, **61**, 208–220.
- 771 42 P. L. Dold, G. A. Ekama, G. Marais and G. R. Marais Van, A general model for the activated
772 sludge process, *Prog. Water Technol.*, 1980, **12**, 47–77.
- 773 43 K. V. Gernaey, M. C. M. Van Loosdrecht, M. Henze, M. Lind and S. B. Jørgensen,
774 Activated sludge wastewater treatment plant modelling and simulation: State of the art,
775 *Environ. Model. Softw.*, 2004, **19**, 763–783.
- 776 44 J. H. Sherrard, R. O. Mines, J. E. Alleman and M. S. Kennedy, Activated Sludge, *J. Water*
777 *Pollut. Control Fed.*, 1983, **55**, 615–622.
- 778 45 P. J. Roeleveld and M. C. M. Van Loosdrecht, Experience with guidelines for wastewater
779 characterisation in The Netherlands | Water Science & Technology, *Water Sci. Technol.*,
780 2002, **45**, 77–87.
- 781 46 M. C. M. Van Loosdrecht, C. M. Lopez-Vazquez, S. C. F. Meijer, C. M. Hooijmans and D.
782 Brdjanovic, Twenty-five years of ASM1: Past, present and future of wastewater treatment
783 modelling, *J. Hydroinformatics*, 2015, **17**, 697–718.
- 784 47 N. Boontian, Cranfield University, 2012.
- 785 48 Z. R. Hu, M. C. Wentzel and G. A. Ekama, Modelling biological nutrient removal activated
786 sludge systems - A review, *Water Res.*, 2003, **37**, 3430–3444.
- 787 49 G. Olsson and B. Newell, Wastewater Treatment Systems: Modelling, Diagnosis and Control.
788 2005, **4**.
- 789 50 A. Serdarevic and A. Dzibur, Wastewater process modeling, *Coupled Syst. Mech.*, 2016, **5**,
790 21–39.
- 791 51 A. G. Dorofeev, Yu. A. Nikolaev, M. N. Kozlov, M. V. Kevbrina, A. M. Agarev, A. Yu.
792 Kallistova and N. V. Pimenov, Modeling of anammox process with the biowin software suite,
793 *Appl. Biochem. Microbiol.*, 2017, **53**, 88–95.
- 794 52 O. Oleyiblo, J. Cao, Q. Feng, G. Wang, X. Zhaoxia and F. Fang, Evaluation and
795 improvement of wastewater treatment plant performance using BioWin, *Chin. J. Oceanol.*
796 *Limnol.*, 2014, **33**, 468–476.
- 797 53 A. Elawwad, M. Matta, M. Abo-Zaid and H. Abdel-Halim, Plant-wide modeling and
798 optimization of a large-scale WWTP using BioWin's ASDM model, *J. Water Process Eng.*,
799 2019, **31**, 100819.
- 800 54 K. Rathore, University of South Florida, 2018.
- 801 55 C. Moragaspiya, J. Rajapakse, W. Senadeera and I. Ali, Simulation of Dynamic Behaviour
802 of a Biological Wastewater Treatment Plant in South East Queensland, Australia using Bio-
803 Win Software, *Eng. J.*, 2017, **21**, 1–22.
- 804 56 R. Vitanza, I. Colussi, A. Cortesi and V. Gallo, Implementing a respirometry-based model
805 into BioWin software to simulate wastewater treatment plant operations, *J. Water Process*
806 *Eng.*, 2016, **9**, 267–275.
- 807 57 I. Hamawand and C. Baillie, Anaerobic Digestion and Biogas Potential: Simulation of Lab
808 and Industrial-Scale Processes, *Energies*, 2015, **8**, 454–474.

- 809 58 J. L. Callahan, Text, Colorado School of Mines, 2018.
- 810 59 K. Venkiteswaran, B. Bocher, J. Maki and D. Zitomer, Relating Anaerobic Digestion
811 Microbial Community and Process Function : Supplementary Issue: Water Microbiology,
812 *Microbiol. Insights*, 2015, **8s2**, MBI.S33593.
- 813 60 J. F. Andrews, A mathematical model for the continuous culture of microorganisms utilizing
814 inhibitory substrates, *Biotechnol. Bioeng.*, 1968, **10**, 707–723.
- 815 61 J. F. Andrews and S. P. Graef, in *Anaerobic Biological Treatment Processes*, American
816 Chemical Society, 1971, pp. 126–162.
- 817 62 D. T. Hill and C. L. Barth, A Dynamic Model for Simulation of Animal Waste Digestion on
818 JSTOR, *Water Pollut. Control Fed.*, 1977, **49**, 2129–2143.
- 819 63 P. Satpathy, S. Steinigeweg, F. Uhlenhut and E. Siefert, Application of Anaerobic Digestion
820 Model 1 (ADM1) for Prediction of Biogas Production, *Int. J. Sci. Eng. Res.*, 2013, **4**, 86–89.
- 821 64 D. J. Batstone, M. Torrijos, C. Ruiz and J. E. Schmidt, Use of an anaerobic sequencing batch
822 reactor for parameter estimation in modelling of anaerobic digestion, *Water Sci. Technol.*,
823 2004, **50**, 295–303.
- 824 65 B. Fezzani and R. Ben Cheikh, Extension of the anaerobic digestion model No. 1 (ADM1) to
825 include phenolic compounds biodegradation processes for the simulation of anaerobic co-
826 digestion of olive mill wastes at thermophilic temperature, *J. Hazard. Mater.*, 2009, **162**,
827 1563–1570.
- 828 66 U. G. Ozkan-Yucel and C. F. Gökçay, Application of ADM1 model to a full-scale anaerobic
829 digester under dynamic organic loading conditions, *Environ. Technol.*, 2010, **31**, 633–640.
- 830 67 A. Ramachandran, R. Rustum and A. J. Adeloje, Review of anaerobic digestion modeling
831 and optimization using nature-inspired techniques, *Processes*, 2019, **7**, 1–12.
- 832 68 J. Palatsi, J. Illa, F. X. Prenafeta-Boldú, M. Laureni, B. Fernandez, I. Angelidaki and X.
833 Flotats, Long-chain fatty acids inhibition and adaptation process in anaerobic thermophilic
834 digestion: Batch tests, microbial community structure and mathematical modelling,
835 *Bioresour. Technol.*, 2010, **101**, 2243–2251.
- 836 69 D. J. Batstone, M. Torrijos, C. Ruiz and J. E. Schmidt, Use of an anaerobic sequencing batch
837 reactor for parameter estimation in modelling of anaerobic digestion, *Water Sci. Technol.*,
838 2004, **50**, 295–303.
- 839 70 E. Shi, J. Li and M. Zhang, Application of IWA Anaerobic Digestion Model No. 1 to
840 simulate butyric acid, propionic acid, mixed acid, and ethanol type fermentative systems
841 using a variable acidogenic stoichiometric approach, *Water Res.*, 2019, **161**, 242–250.
- 842 71 A. E. Rotaru, P. M. Shrestha, F. Liu, M. Shrestha, D. Shrestha, M. Embree, K. Zengler, C.
843 Wardman, K. P. Nevin and D. R. Lovley, A new model for electron flow during anaerobic
844 digestion: Direct interspecies electron transfer to *Methanosaeta* for the reduction of carbon
845 dioxide to methane, *Energy Environ. Sci.*, 2014, **7**, 408–415.
- 846 72 D. R. Lovley, Syntrophy Goes Electric: Direct Interspecies Electron Transfer, *Annu. Rev.*
847 *Microbiol.*, 2017, **71**, 643–664.
- 848 73 T. Storck, B. Virdis and D. J. Batstone, Modelling extracellular limitations for mediated
849 versus direct interspecies electron transfer, *ISME J.*, 2016, **10**, 621–631.
- 850 74 Y. Liu, Y. Zhang, Z. Zhao, H. H. Ngo, W. Guo, J. Zhou, L. Peng and B. J. Ni, A modeling
851 approach to direct interspecies electron transfer process in anaerobic transformation of
852 ethanol to methane, *Environ. Sci. Pollut. Res.*, 2017, **24**, 855–863.
- 853 75 A. Terada, S. Zhou and M. Hosomi, Presence and detection of anaerobic ammonium-
854 oxidizing (anammox) bacteria and appraisal of anammox process for high-strength

- 855 nitrogenous wastewater treatment: a review, *Clean Technol. Environ. Policy*, 2011, **13**, 759–
856 781.
- 857 76 X. Chang, D. Li, Y. Liang, Z. Yang, S. Cui, T. Liu, H. Zeng and J. Zhang, Performance of a
858 completely autotrophic nitrogen removal over nitrite process for treating wastewater with
859 different substrates at ambient temperature, *J. Environ. Sci.*, 2013, **25**, 688–697.
- 860 77 K. A. Third, A. O. Sliekers, J. G. Kuenen and M. S. M. Jetten, The CANON system
861 (completely autotrophic nitrogen-removal over nitrite) under ammonium limitation:
862 Interaction and competition between three groups of bacteria, *Syst. Appl. Microbiol.*, 2001,
863 **24**, 588–596.
- 864 78 G. Cema, A. Sochacki, J. Kubiawicz, P. Gutwiński and J. Surmacz-Górska, Start-up,
865 modelling and simulation of the anammox process in a membrane bioreactor, *Chem. Process
866 Eng. - Inzynieria Chem. Proces.*, 2012, **33**, 639–650.
- 867 79 B.-J. Ni, A. Joss and Z. Yuan, Modeling nitrogen removal with partial nitrification and
868 anammox in one floc-based sequencing batch reactor, 2014, **67**, 321–329.
- 869 80 D. Pant, A. Singh, G. Van Bogaert, S. Irving Olsen, P. Singh Nigam, L. Diels and K.
870 Vanbroekhoven, *RSC Adv.*, 2012, **2**, 1248–1263.
- 871 81 C. Xia, D. Zhang, W. Pedrycz, Y. Zhu and Y. Guo, *J. Power Sources*, 2018, 373, 119–131.
- 872 82 V. B. Oliveira, M. Simões, L. F. Melo and A. M. F. R. Pinto, Overview on the developments
873 of microbial fuel cells, *Biochem. Eng. J.*, 2013, **73**, 53–64.
- 874 83 V. M. Ortiz-Martínez, M. J. Salar-García, A. P. de los Ríos, F. J. Hernández-Fernández,
875 J. A. Egeac and L. J. Lozano, Developments in microbial fuel cell modeling, *Chem. Eng. J.*,
876 2015, **271**, 50–60.
- 877 84 R. P. Pinto, B. Srinivasan, A. Escapa and B. Tartakovsky, Multi-population model of a
878 microbial electrolysis cell, *Environ. Sci. Technol.*, 2011, **45**, 5039–5046.
- 879 85 R. P. Pinto, B. Srinivasan, M. F. Manuel and B. Tartakovsky, A two-population bio-
880 electrochemical model of a microbial fuel cell, *Bioresour. Technol.*, 2010, **101**, 5256–5265.
- 881 86 Q. Ping, C. Zhang, X. Chen, B. Zhang, Z. Huang and Z. He, Mathematical Model of
882 Dynamic Behavior of Microbial Desalination Cells for Simultaneous Wastewater Treatment
883 and Water Desalination, *Environ. Sci. Technol.*, 2014, **48**, 13010–13019.
- 884 87 A. K. Marcus, C. I. Torres and B. E. Rittmann, Conduction-based modeling of the biofilm
885 anode of a microbial fuel cell, *Biotechnol. Bioeng.*, 2007, **98**, 1171–1182.
- 886 88 F. Harnisch, R. Warmbier, R. Schneider and U. Schröder, Modeling the ion transfer and
887 polarization of ion exchange membranes in bioelectrochemical systems, *Bioelectrochemistry*,
888 2009, **75**, 136–141.
- 889 89 J. Kim, H. Kim, B. Kim and J. Yu, Computational fluid dynamics analysis in microbial fuel
890 cells with different anode configurations, *Water Sci. Technol.*, 2014, **69**, 1447–1452.
- 891 90 D. C. Vuono, J. Regnery, D. Li, Z. L. Jones, R. W. Holloway and J. E. Drewes, rRNA Gene
892 Expression of Abundant and Rare Activated-Sludge Microorganisms and Growth Rate
893 Induced Micropollutant Removal, *Environ. Sci. Technol.*, 2016, **50**, 6299–6309.
- 894 91 L. Wu, D. Ning, B. Zhang, Y. Li, P. Zhang, X. Shan, Q. Zhang, M. Brown, Z. Li, J. D. Van
895 Nostrand, F. Ling, N. Xiao, Y. Zhang, J. Vierheilig, G. F. Wells, Y. Yang, Y. Deng, Q. Tu,
896 A. Wang, T. Zhang, Z. He, J. Keller, P. H. Nielsen, P. J. J. Alvarez, C. S. Criddle, M.
897 Wagner, J. M. Tiedje, Q. He, T. P. Curtis, D. A. Stahl, L. Alvarez-Cohen, B. E. Rittmann, X.
898 Wen, J. Zhou, D. Acevedo, M. Agullo-Barcelo, G. L. Andersen, J. C. de Araujo, K. Boehnke,
899 P. Bond, C. B. Bott, P. Bovio, R. K. Brewster, F. Bux, A. Cabezas, L. Cabrol, S. Chen, C.
900 Etchebehere, A. Ford, D. Frigon, J. S. GÃmez, J. S. Griffin, A. Z. Gu, M. Habagil, L. Hale,

- 901 S. D. Hardeman, M. Harmon, H. Horn, Z. Hu, S. Jauffur, D. R. Johnson, A. Keucken, S.
902 Kumari, C. D. Leal, L. A. Lebrun, J. Lee, M. Lee, Z. M. P. Lee, M. Li, X. Li, Y. Liu, R. G.
903 Luthy, L. C. Mendonça-Sa-Hagler, F. G. R. de Menezes, A. J. Meyers, A. Mohebbi, A.
904 Oehmen, A. Palmer, P. Parameswaran, J. Park, D. Patsch, V. Reginatto, F. L. de los Reyes,
905 A. Noyola, S. Rossetti, J. Sidhu, W. T. Sloan, K. Smith, O. V. de Sousa, K. Stephens, R.
906 Tian, N. B. Tooker, D. De los Cobos Vasconcelos, S. Wakelin, B. Wang, J. E. Weaver, S.
907 West, P. Wilmes, S. G. Woo, J. H. Wu, L. Wu, C. Xi, M. Xu, T. Yan, M. Yang, M. Young,
908 H. Yue, Q. Zhang, W. Zhang, Y. Zhang and H. Zhou, Global diversity and biogeography of
909 bacterial communities in wastewater treatment plants, *Nat. Microbiol.*, 2019, **4**, 1183–1195.
- 910 92 A. M. Saunders, M. Albertsen, J. Vollertsen and P. H. Nielsen, The activated sludge
911 ecosystem contains a core community of abundant organisms, *ISME J.*, 2016, **10**, 11–20.
- 912 93 S. Jung and J. M. Regan, Comparison of anode bacterial communities and performance in
913 microbial fuel cells with different electron donors, *Appl. Microbiol. Biotechnol.*, 2007, **77**,
914 393–402.
- 915 94 D. Pant, G. Van Bogaert, L. Diels and K. Vanbroekhoven, A review of the substrates used in
916 microbial fuel cells (MFCs) for sustainable energy production, *Bioresour. Technol.*, 2010,
917 **101**, 1533–1543.
- 918 95 S. Ishii, S. Suzuki, T. M. Norden-Krichmar, A. Tenney, P. S. G. Chain, M. B. Scholz, K. H.
919 Neilson and O. Bretschger, A novel metatranscriptomic approach to identify gene
920 expression dynamics during extracellular electron transfer, *Nat. Commun.*, 2013, **4**, 1–10.
- 921 96 B. E. Logan, Exoelectrogenic bacteria that power microbial fuel cells, *Nat. Rev. Microbiol.*,
922 2009, **7**, 375–381.
- 923 97 Y. Xiao, Y. Zheng, S. Wu, E.-H. Zhang, Z. Chen, P. Liang, X. Huang, Z.-H. Yang, I.-S. Ng,
924 B.-Y. Chen and F. Zhao, Pyrosequencing Reveals a Core Community of Anodic Bacterial
925 Biofilms in Bioelectrochemical Systems from China, *Front. Microbiol.*, 2015, **6**, 1410.
- 926 98 X. Zhao, L. Li, D. Wu, T. Xiao, Y. Ma and X. Peng, Modified Anaerobic Digestion Model
927 No. 1 for modeling methane production from food waste in batch and semi-continuous
928 anaerobic digestions, *Bioresour. Technol.*, 2019, **271**, 109–117.
- 929 99 P. Reichert and N. Schuwirth, Linking statistical bias description to multiobjective model
930 calibration, *Water Resour. Res.*, 2012, **48**, 2011WR011391.
- 931 100 X. Flotats, B. K. Ahring and I. Angelidaki, Parameter identification of thermophilic
932 anaerobic degradation of valerate, *Appl. Biochem. Biotechnol.*, 2003, **109**, 47–62.
- 933 101 N. Noykova, T. Müller, M. Gyllenberg and J. Timmer, Quantitative analyses of anaerobic
934 wastewater treatment processes: identifiability and parameter estimation., *Biotechnol.*
935 *Bioeng.*, 2002, **78**, 89–103.
- 936 102 A. Donoso-Bravo, J. Mailier, C. Martin, J. Rodríguez, C. A. Aceves-Lara and A. V. Wouwer,
937 Model selection, identification and validation in anaerobic digestion: A review, *Water Res.*,
938 2011, **45**, 5347–5364.
- 939 103 T. G. Müller, N. Noykova, M. Gyllenberg and J. Timmer, Parameter identification in
940 dynamical models of anaerobic waste water treatment, *Math. Biosci.*, 2002, **177–178**, 147–
941 160.
- 942 104 Z. Boger, Application of neural networks to water and wastewater treatment plant operation,
943 *ISA Trans.*, 1992, **31**, 25–33.
- 944 105 J. Jawad, A. H. Hawari and S. Javaid Zaidi, Artificial neural network modeling of
945 wastewater treatment and desalination using membrane processes: A review, *Chem. Eng. J.*,
946 2021, **419**, DOI:10.1016/j.cej.2021.129540.

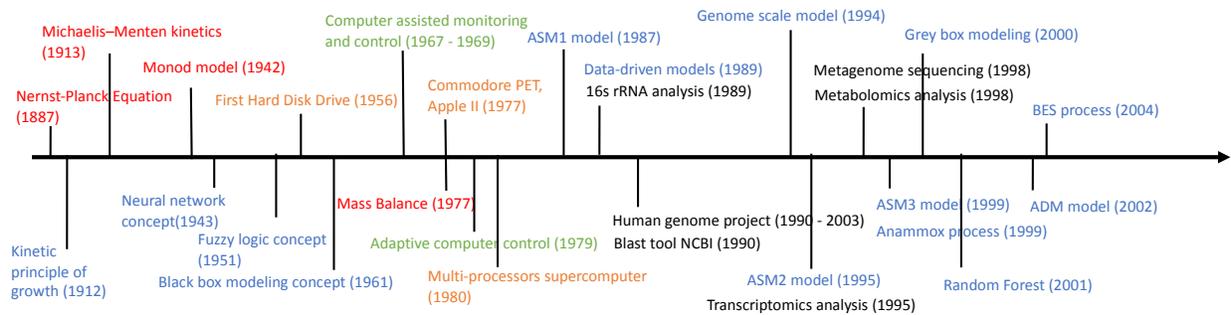
- 947 106S. Zendejboudi, N. Rezaei and A. Lohi, *Appl. Energy*, 2018, 228, 2539–2566.
- 948 107M. Ben Nasr and M. Chtourou, A hybrid training algorithm for feedforward neural networks,
949 *Neural Process. Lett.*, 2006, **24**, 107–117.
- 950 108J. Thibault, V. Van Breusegem and A. Chéruey, On-line Prediction of Fermentation Variables
951 Using Neural Networks, *Biotechnol. Bioeng.*, 1990, **36**, 1041–1048.
- 952 109Y. Y. Yang and D. A. Linkenst, Modelling of continuous bioreactors via neural networks,
953 *Trans. Inst. Meas. Control*, 1993, **15**, 158–169.
- 954 110P. Kundu, A. Debsarkar and S. Mukherjee, Artificial Neural Network Modeling for
955 Biological Removal of Organic Carbon and Nitrogen from Slaughterhouse Wastewater in a
956 Sequencing Batch Reactor, *Adv. Artif. Neural Syst.*, 2013, **2013**, 1-15,
957 DOI:10.1155/2013/268064.
- 958 111B. Mahanty, M. Zafar and H. S. Park, Characterization of co-digestion of industrial sludges
959 for biogas production by artificial neural network and statistical regression models, *Environ.*
960 *Technol. U. K.*, 2013, **34**, 2145–2153.
- 961 112P. Koehn, University of Tennessee, 1994.
- 962 113M. Shariati, M. S. Mafipour, P. Mehrabi, A. Bahadori, Y. Zandi, M. N. A. Salih, H. Nguyen,
963 J. Dou, X. Song and S. Poi-Ngian, Application of a hybrid artificial neural network-particle
964 swarm optimization (ANN-PSO) model in behavior prediction of channel shear connectors
965 embedded in normal and high-strength concrete, *Appl. Sci. Switz.*, 2019, **9**,
966 DOI:10.3390/app9245534.
- 967 114H. Lu, J. Chen and L. Guo, in *Comprehensive Energy Systems*, Elsevier, 2018, vol. 5–5, pp.
968 258–314.
- 969 115H. Tyrallis, G. Papacharalampous and A. Langousis, A brief review of random forests for
970 water scientists and practitioners and their recent history in water resources, *Water*, 2019, **11**,
971 910.
- 972 116L. Breiman, *RANDOM FORESTS-RANDOM FEATURES*, 1999.
- 973 117W. Li, C. Li and T. Wang, Application of machine learning algorithms in mbr simulation
974 under big data platform, *Water Pract. Technol.*, 2020, **15**, 1238–1247.
- 975 118M. J. Song, S. Choi, W. Bin Bae, J. Lee, H. Han, D. D. Kim, M. Kwon, J. Myung, Y. M.
976 Kim and S. Yoon, Identification of primary effecters of N₂O emissions from full-scale
977 biological nitrogen removal systems using random forest approach, *Water Res.*, 2020, 184,
978 116114, DOI:10.1016/j.watres.2020.116144.
- 979 119B. Szeląg, A. Gawdzik and A. Gawdzik, Application of selected methods of black box for
980 modeling the settlability process in wastewater treatment plant, 2017, **24**, 119-127,
981 DOI:10.1515/eces-2017-0009.
- 982 120B. Szeląg, J. Drewnowski, G. Łagód, D. Majerek, E. Dacewicz and F. Fatone, Soft sensor
983 application in identification of the activated sludge bulking considering the technological and
984 economical aspects of smart systems functioning, *Sens. Switz.*, 2020, **20**, 1941.
- 985 121H. Tyrallis, G. Papacharalampous and A. Langousis, A Brief Review of Random Forests for
986 Water Scientists and Practitioners and Their Recent History in Water Resources, *Water*,
987 2019, **11**, 910.
- 988 122M. J. Song, S. Choi, W. B. Bae, J. Lee, H. Han, D. D. Kim, M. Kwon, J. Myung, Y. M. Kim
989 and S. Yoon, Identification of primary effecters of N₂O emissions from full-scale biological
990 nitrogen removal systems using random forest approach, *Water Res.*, 2020, **184**, 116114.

- 991 123D. Wang, S. Thunéll, U. Lindberg, L. Jiang, J. Trygg, M. Tysklind and N. Souihi, A machine
992 learning framework to improve effluent quality control in wastewater treatment plants, *Sci.*
993 *Total Environ.*, 2021, **784**, 147138.
- 994 124G. Louppe, University of Liège, 2014.
- 995 125J. M. Mendel and G. C. Mouzouris, Designing Fuzzy Logic Systems, *IEEE Trans. Circuits*
996 *Syst.*, 1997, **44**, 885.
- 997 126L. Fan and K. Boshnakov, in *Proceedings of the World Congress on Intelligent Control and*
998 *Automation (WCICA)*, 2010, pp. 4142–4146.
- 999 127H. Xu and R. Vilanova, in *2015 23rd Mediterranean Conference on Control and Automation,*
1000 *MED 2015 - Conference Proceedings*, Institute of Electrical and Electronics Engineers Inc.,
1001 2015, pp. 545–550.
- 1002 128T.J.J.Kalker, C.P.van Goor, P.J.Roeleveld, M.F.Ruland and R.Babuška, Fuzzy control of
1003 aeration in an activated sludge wastewater treatment plant: design, simulation and evaluation,
1004 *Water Sci. Technol.*, 1999, **39**, 61–69.
- 1005 129A. Robles, E. Latrille, M. V. Ruano and J. P. Steyer, A fuzzy-logic-based controller for
1006 methane production in anaerobic fixed-film reactors, *Public Health Titles*, 2017, **38**, 42–52.
- 1007 130J. S. R. Jang, ANFIS: Adaptive-Network-Based Fuzzy Inference System, *IEEE Trans. Syst.*
1008 *Man Cybern.*, 1993, **23**, 665–685.
- 1009 131R. Rustum, *Modelling Activated Sludge Wastewater Treatment Plants Using Artificial*
1010 *Intelligence Techniques (Fuzzy Logic and Neural Networks)*, 2009.
- 1011 132T. Y. Pai, P. Y. Yang, S. C. Wang, M. H. Lo, C. F. Chiang, J. L. Kuo, H. H. Chu, H. C. Su, L.
1012 F. Yu, H. C. Hu and Y. H. Chang, Predicting effluent from the wastewater treatment plant of
1013 industrial park based on fuzzy network and influent quality, *Appl. Math. Model.*, 2011, **35**,
1014 3674–3684.
- 1015 133M. Huang, J. Wan, K. Hu, Y. Ma and Y. Wang, Enhancing dissolved oxygen control using
1016 an on-line hybrid fuzzy-neural soft-sensing model-based control system in an
1017 anaerobic/anoxic/oxic process, *J Ind Microbiol Biotechnol*, 2013, **40**, 1393–1401.
- 1018 134I. A. Essienubong, A.-I. Effiong Ndon and J. Etim, Fuzzy modeling and optimization of
1019 anaerobic co-digestion process parameters for effective biogas yield from bio-wastes, *Int. J.*
1020 *Energy Eng. Sci.*, 2020, 43–61.
- 1021 135A. Hosseinzadeh, J. L. Zhou, A. Altaee, M. Baziar and X. Li, Modeling water flux in
1022 osmotic membrane bioreactor by adaptive network-based fuzzy inference system and
1023 artificial neural network, *Bioresour. Technol.*, 2020, **310**, 123391,
1024 DOI:10.1016/j.biortech.2020.123391.
- 1025 136M. Ansari, F. Othman and A. El-Shafie, Optimized fuzzy inference system to enhance
1026 prediction accuracy for influent characteristics of a sewage treatment plant, *Sci. Total*
1027 *Environ.*, 2020, **722**, 137878.
- 1028 137K. Yetilmezsoy, in *Handbook of Environmental Materials Management*, Springer
1029 International Publishing, 2019, pp. 2001–2046.
- 1030 138C. González-Figueroa, R. Alejandro Flores-Estrella and O. A. Rojas-Rejón, in *Current*
1031 *Topics in Biochemical Engineering*, IntechOpen, 2019.
- 1032 139Y. D. Zhao, E. N. Yamoah and P. G. Gillespie, Regeneration of broken tip links and
1033 restoration of mechanical transduction in hair cells, *Proc. Natl. Acad. Sci. U. S. A.*, 1996, **93**,
1034 15469–15474.
- 1035 140D. S. Lee, C. O. Jeon, J. M. Park and K. S. Chang, Hybrid neural network modeling of a full-
1036 scale industrial wastewater treatment process, *Biotechnol. Bioeng.*, 2002, **78**, 670–682.

- 1037 141 D. S. Lee, P. A. Vanrolleghem and M. P. Jong, Parallel hybrid modeling methods for a full-
1038 scale cokes wastewater treatment plant, *J. Biotechnol.*, 2005, **115**, 317–328.
- 1039 142 M. von Stosch, R. Oliveira, J. Peres and S. Feyo de Azevedo, Hybrid semi-parametric
1040 modeling in process systems engineering: Past, present and future, *Comput. Chem. Eng.*,
1041 2014, **60**, 86–101.
- 1042 143 S. Banihashemi, G. Ding and J. Wang, Developing a Hybrid Model of Prediction and
1043 Classification Algorithms for Building Energy Consumption, *Energy Procedia*, 2017, **110**,
1044 371–376.
- 1045 144 S. J. McIlroy, A. M. Saunders, M. Albertsen, M. Nierychlo, B. McIlroy, A. A. Hansen, S. M.
1046 Karst, J. L. Nielsen and P. H. Nielsen, MiDAS: the field guide to the microbes of activated
1047 sludge, *Database*.
- 1048 145 N. Hvala and J. Kocijan, Design of a hybrid mechanistic/Gaussian process model to predict
1049 full-scale wastewater treatment plant effluent, *Comput. Chem. Eng.*, 2020, **140**, 106914,
1050 DOI:10.1016/j.compchemeng.2020.106934.
- 1051 146 T. Větrovský and P. Baldrian, The Variability of the 16S rRNA Gene in Bacterial Genomes
1052 and Its Consequences for Bacterial Community Analyses, *PLoS ONE*, 2013, **8**, 1–10,
1053 DOI:10.1371/journal.pone.0057923.
- 1054 147 V. Aguiar-Pulido, W. Huang, V. Suarez-Ulloa, T. Cickovski, K. Mathee and G. Narasimhan,
1055 *Evol. Bioinforma.*, 2016, **12**, 5–16.
- 1056 148 K. Yu and T. Zhang, Metagenomic and metatranscriptomic analysis of microbial community
1057 structure and gene expression of activated sludge, *PLoS ONE*, 2012, **7**, 38183.
- 1058 149 L. Ye, R. Mei, W. T. Liu, H. Ren and X. X. Zhang, Machine learning-aided analyses of
1059 thousands of draft genomes reveal specific features of activated sludge processes,
1060 *Microbiome*, 2020, **8**, 1–13.
- 1061 150 P. E. Larsen, D. Field and J. A. Gilbert, Predicting bacterial community assemblages using
1062 an artificial neural network approach, *Nat. Methods*, 2012, **9**, 621–625.
- 1063 151 J. L. Metcalf, Z. Z. Xu, S. Weiss, S. Lax, W. Van Treuren, E. R. Hyde, S. J. Song, A. Amir,
1064 P. Larsen, N. Sangwan, D. Haarmann, G. C. Humphrey, G. Ackermann, L. R. Thompson, C.
1065 Lauber, A. Bibat, C. Nicholas, M. J. Gebert, J. F. Petrosino, S. C. Reed, J. A. Gilbert, A. M.
1066 Lynne, S. R. Bucheli, D. O. Carter and R. Knight, Microbial community assembly and
1067 metabolic function during mammalian corpse decomposition, *Science*, 2016, **351**, 158–162.
- 1068 152 K. L. Lesnik and H. Liu, Predicting Microbial Fuel Cell Biofilm Communities and
1069 Bioreactor Performance using Artificial Neural Networks, *Environ. Sci. Technol.*, 2017, **51**,
1070 10881–10892.
- 1071 153 K. L. Lesnik, W. Cai and H. Liu, Microbial Community Predicts Functional Stability of
1072 Microbial Fuel Cells, *Environ. Sci. Technol.*, , DOI:10.1021/acs.est.9b03667.
- 1073 154 H. Yuan, S. Sun, I. M. Abu-Reesh, B. D. Badgley and Z. He, Unravelling and
1074 Reconstructing the Nexus of Salinity, Electricity, and Microbial Ecology for
1075 Bioelectrochemical Desalination, *Environ. Sci. Technol.*, 2017, **51**, 12672–12682.
- 1076 155 J. Kuang, L. Huang, Z. He, L. Chen, Z. Hua, P. Jia, S. Li, J. Liu, J. Li, J. Zhou and W. Shu,
1077 Predicting taxonomic and functional structure of microbial communities in acid mine
1078 drainage, *ISME J.*, 2016, **10**, 1527–1539.
- 1079 156 H. Yuan, R. Mei, J. Liao and W. T. Liu, Nexus of stochastic and deterministic processes on
1080 microbial community assembly in biological systems, *Front. Microbiol.*, 2019, **10**, 1–12,
1081 DOI:10.3389/fmicb.2019.01536.

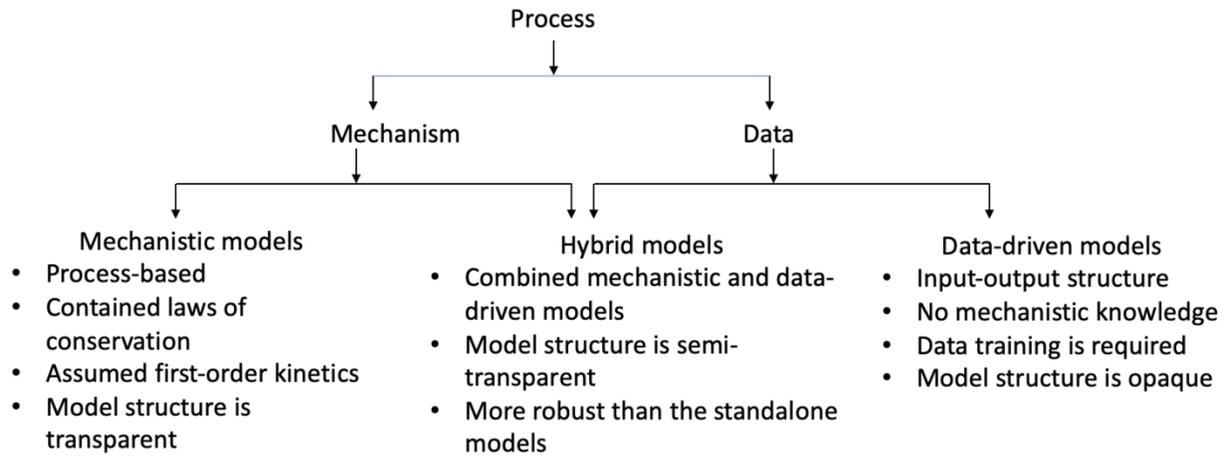
- 1082 157I. Vanwonterghem, P. D. Jensen, P. G. Dennis, P. Hugenholtz, K. Rabaey and G. W. Tyson,
 1083 Deterministic processes guide long-term synchronised population dynamics in replicate
 1084 anaerobic digesters, *ISME J.*, 2014, **8**, 2015–2028.
 1085 158J. S. Griffin and G. F. Wells, Regional synchrony in full-scale activated sludge bioreactors
 1086 due to deterministic microbial community assembly, *ISME J.*, 2017, **11**, 500–511.
 1087 159L. Zhang, P. Zheng, C. Tang and R. Jin, Anaerobic ammonium oxidation for treatment of
 1088 ammonium-rich wastewaters., *J. Zhejiang Univ. Sci. B*, 2008, **9**, 416–426.
 1089

1090



1091 Figure 1. Timeline of model concepts (red), model strategies (blue), control and monitoring
 1092 development (green), computational infrastructure (orange), and development of molecular
 1093 biology and bioinformatics (black).
 1094

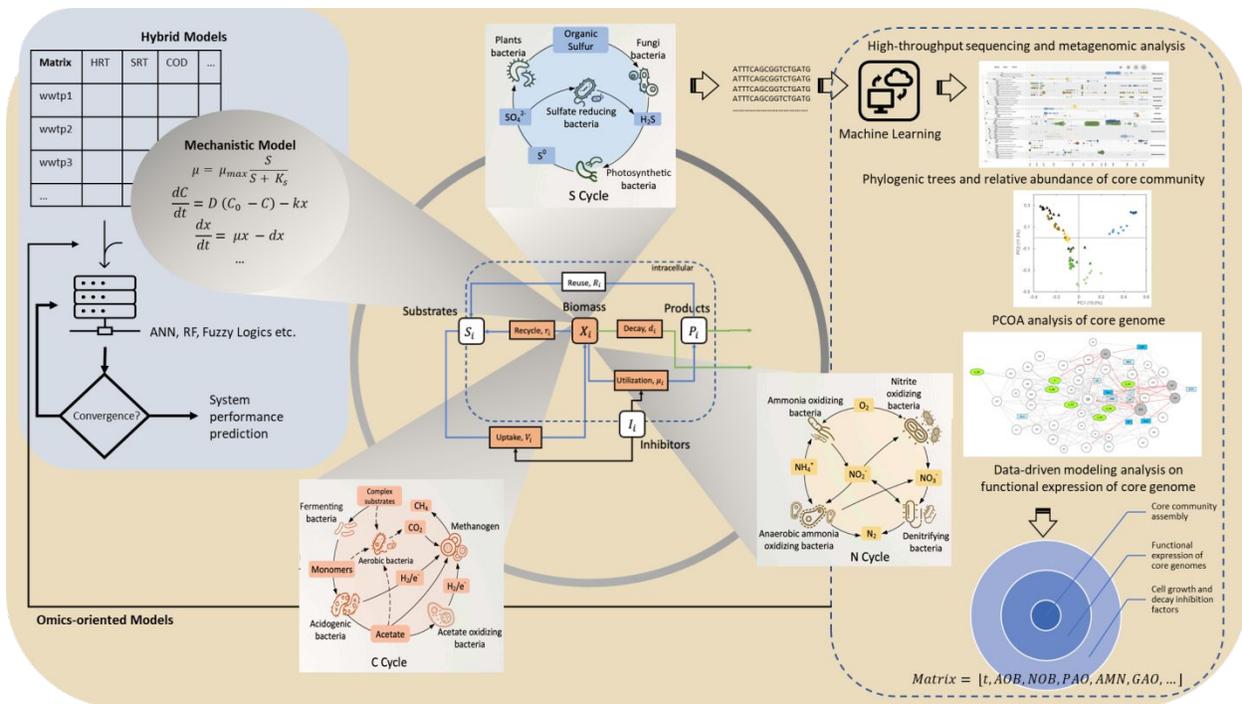
1095



1096

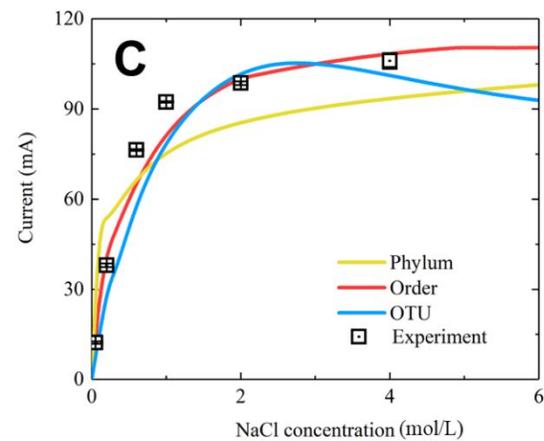
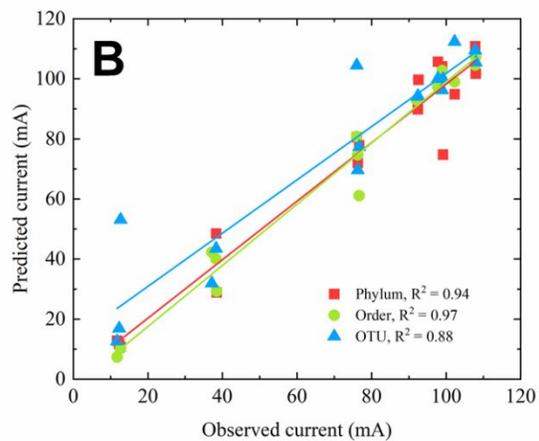
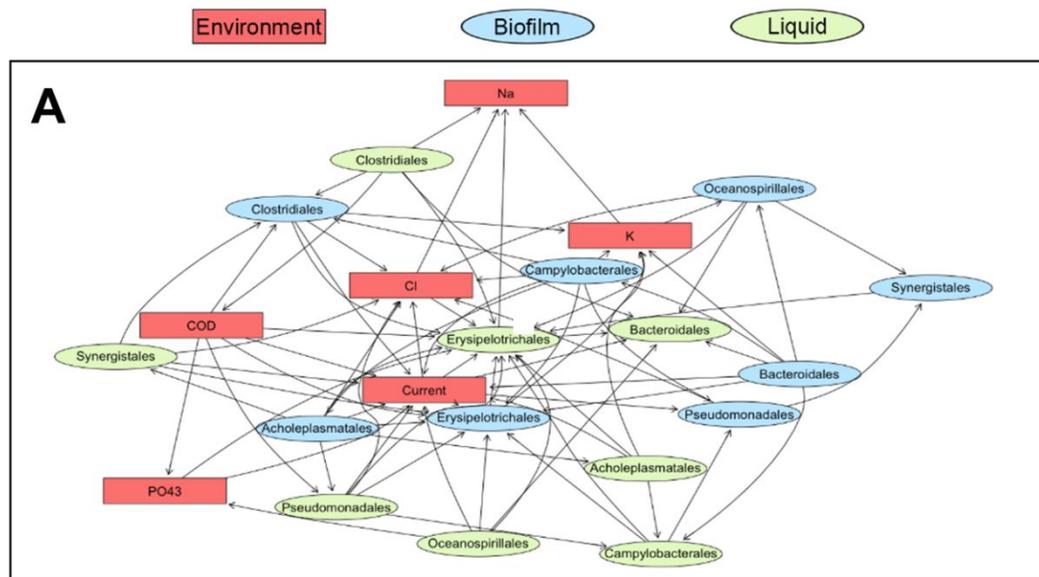
1097 Figure 2. The characteristics and advantages of mechanistic, data-driven, and hybrid models

1098



1099

1100 Figure 3. Incorporation of genomic data into model construction. The diagrams of S-, N-, and C-
 1101 cycle are originated from the study of Wu and Yin⁷. The figures of phylogenetic trees, PCoA, and
 1102 data-driven modeling analysis on functional expression are based on the study of Cheng et al.³¹



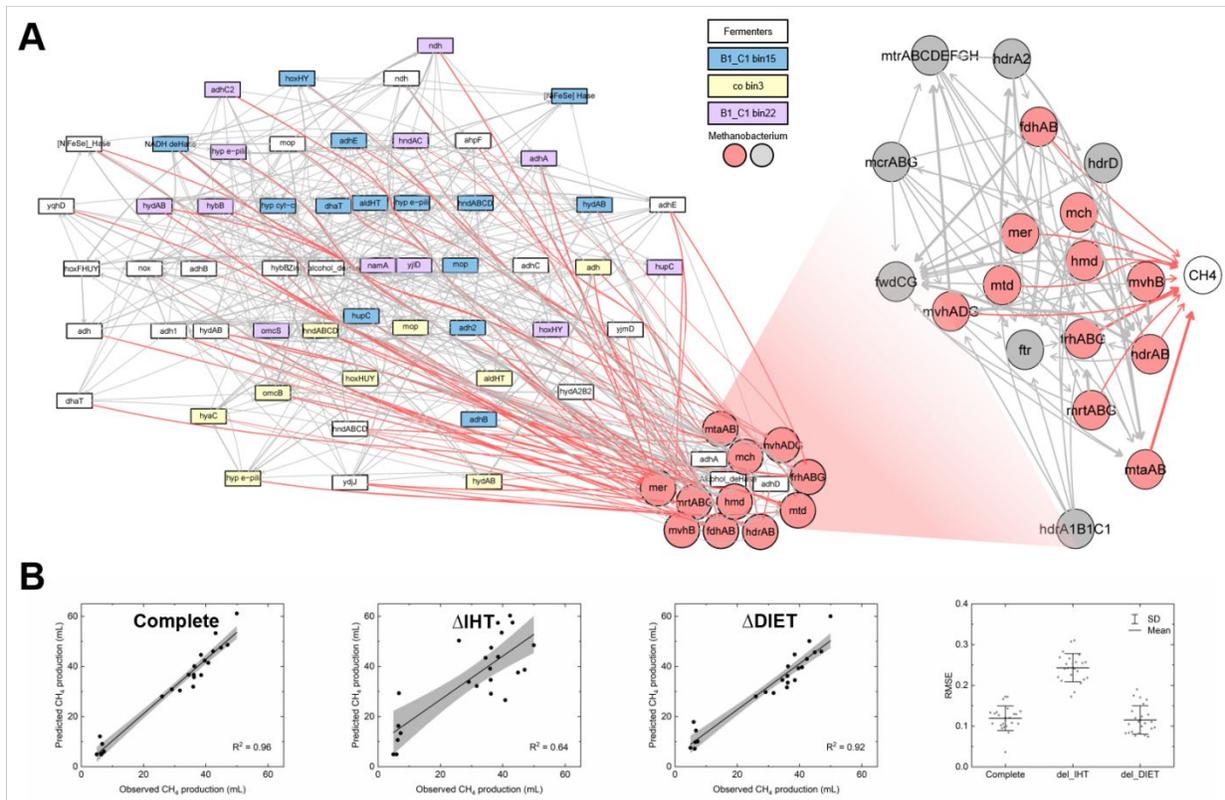
1103

1104 Figure 4. (A) A Bayesian network trained with the microbial population dynamics at the order

1105 level in a bioelectrochemical system. (B) Predicted vs. observed current production. (C)

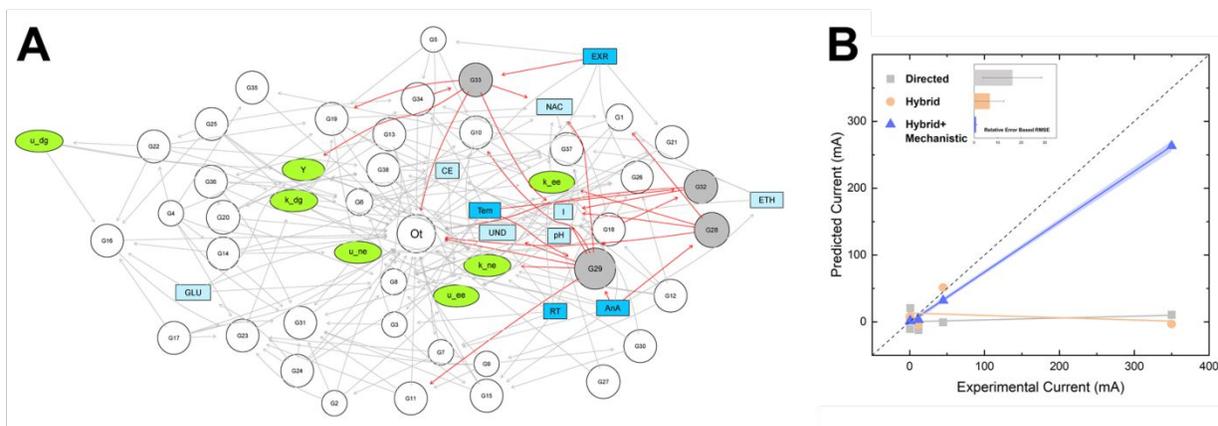
1106 Prediction of current production as a function of substrate salinity. Figures adapted from Yuan et

1107 al.¹⁵⁴



1108

1109 Figure 5. (A) A Bayesian network trained with the genes for alcohol metabolism, hydrogen
 1110 metabolism, direct interspecies electron transfer, and methanogenesis from dominant microbes.
 1111 (B) Prediction of methane production with a complete Bayesian network and in-silico knockout
 1112 of relevant genes. Figures adapted from Yuan et al.³⁰



1113

1114 Figure 6. (A) A Bayesian network as the data-driven component of the hybrid model was trained

1115 with microbial population dynamics and microbial kinetic parameters estimated from the

1116 mechanistic component (green oval nodes). (B) Predicted vs. observed current production.

1117 Figures adapted from Cheng et al.³¹