



**HAL**  
open science

## **Data assimilation for prediction of ammonium in wastewater treatment plant: From physical to data driven models**

Victor Bertret, Roman Le Goff Latimier, Valérie Monbet

### ► **To cite this version:**

Victor Bertret, Roman Le Goff Latimier, Valérie Monbet. Data assimilation for prediction of ammonium in wastewater treatment plant: From physical to data driven models. *Water Research*, 2025, 282, pp.123673. <10.1016/j.watres.2025.123673>. <hal-05055967>

**HAL Id: hal-05055967**

**<https://hal.science/hal-05055967v1>**

Submitted on 10 Jun 2025

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons CC BY 4.0 - Attribution - International License

1 Data assimilation for prediction of ammonium in  
2 wastewater treatment plant: from physical to data  
3 driven models

4 Victor Bertret<sup>a,b,c,\*</sup>, Roman Le Goff Latimier<sup>b</sup>, Valérie Monbet<sup>c</sup>

*<sup>a</sup>Purecontrol, 68 Av. Sergent Maginot, Rennes, 35000, FRANCE*

*<sup>b</sup>SATIE, ENS Rennes, CNRS, Bruz, 35170, FRANCE*

*<sup>c</sup>IRMAR, Université Rennes 1, CNRS, Rennes, 35700, FRANCE*

---

5 **Abstract**

6 This study compares various modeling approaches to predict ammonium  
7 concentration in wastewater treatment plants (WWTPs), with a focus on  
8 integrating data assimilation techniques. It explores white-box, grey-box,  
9 and black-box models, evaluating their ability to capture the complex dy-  
10 namics of WWTPs and manage uncertainties associated with limited data  
11 and sensor noise. The article highlights the importance of data assimilation  
12 for simultaneously calibrating model parameters, latent variables (such as  
13 unmeasured species concentrations), and quantifying prediction uncertainty.  
14 Simulation results demonstrate that the non-parametric black box model  
15 outperforms all other models in terms of predictive accuracy and uncertainty  
16 estimation. This finding underscores the effectiveness of machine learning  
17 when integrated with data assimilation techniques to extract insights from  
18 training datasets, even in the presence of limited data. Interestingly, the  
19 addition of an extra sensor, such as an oxygen sensor, did not enhance model  
20 performance. Experiments conducted in a real system showed that the non-  
21 parametric black box model could effectively capture the general dynamics of

22 ammonium concentration in an actual wastewater treatment plant. However,  
23 its performance was somewhat diminished compared to simulation results,  
24 likely due to variability in input concentrations that were not accounted for  
25 in the model.

26 *Keywords:* Activated Sludge, Data assimilation, Model Identification,  
27 Uncertainty Quantification, Expectation Maximization, Conditional  
28 Particle Filter

---

## 29 **1. Introduction**

30 Wastewater treatment plants (WWTPs) are complex, energy-intensive  
31 systems that must meet increasingly stringent regulations regarding efflu-  
32 ent quality and energy efficiency. Meanwhile, energy prices and volatility  
33 have increased significantly in recent years (RTE, 2021). One of the highest  
34 energy demands at these facilities is the aeration process (Reardon, 1995),  
35 where, among other processes, nitrification takes place. Therefore, a signifi-  
36 cant challenge is to reduce the concentration of ammonium ( $NH_4^+$ ), primarily  
37 originating from the incoming water, while simultaneously minimizing energy  
38 consumption. The Alternating Activated Sludge Plant (AASP) is a specific  
39 configuration of activated sludge plants (ASPs) that achieves complete ni-  
40 trogen removal in a single bioreactor through nitrification and denitrification  
41 (Chachuat et al., 2005), as shown in Figure 1. It alternates between aeration  
42 (introducing oxygen ( $O_2$ ) into the aeration tank with compressors) and non-  
43 aeration phases, requiring precise control to meet strict regulations, including  
44 daily average ammonium ( $NH_4^+$ ) limits (Balku, 2007; Lukasse et al., 1999).

45 Given the complex dynamics and strict regulatory requirements of AASPs,  
46 modeling serves as a valuable approach for developing optimization strate-

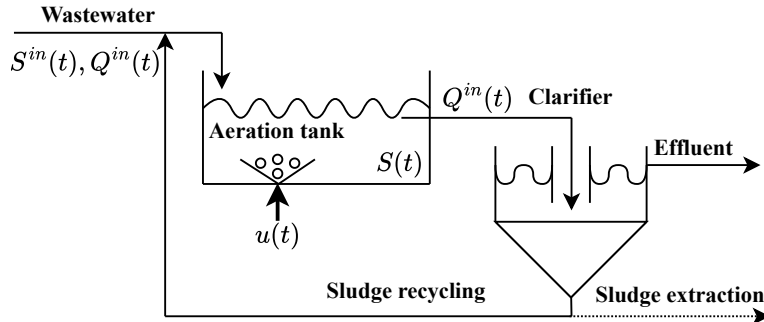


Figure 1: A schematic representation of the secondary treatment process in an AASP plant. The study focuses on quantifying the concentration of the species in the aeration tank, denoted by  $S(t)$ , which is affected by the influent wastewater concentration  $S^{in}(t)$ , the inflow rate  $Q^{in}(t)$ , and the control input  $u(t)$ , representing the state of the compressors.

47 gies. While deep reinforcement learning (DRL) has been explored as a control  
 48 strategy for wastewater treatment plants (WWTPs) (Fu et al., 2022; Sutton  
 49 and Barto, 2018), real-world applications often struggle with stability and  
 50 reliability issues (Dulac-Arnold et al., 2019). In contrast, model-based ap-  
 51 proaches can provide more reliable control by providing a better understand-  
 52 ing of system dynamics. However, these approaches encounter significant  
 53 challenges in AASPs due to the complex interactions within the system and  
 54 variability in influent composition, further complicated by limited and noisy  
 55 sensor data (Lukasse et al., 1999). It is important to distinguish between  
 56 estimating the varying influent composition and reconstructing the system’s  
 57 dynamics. Although influent estimation is well-documented in the literature  
 58 (Andreides et al., 2022), it remains outside the scope of this study. Instead,  
 59 our focus is on the reconstruction of system dynamics in ASPs, which presents  
 60 several challenges:

- 61     • **Medium Dimensionality:** although there are many species to track,  
62       the number of variables is relatively moderate compared to fields like  
63       climate or oceanography.
- 64     • **High Interaction and Uncertainty:** the concentrations of species  
65       interact with each other in complex ways, influenced by unobserved  
66       factors such as pH, temperature, and sludge characteristics.
- 67     • **Partial Observation:** the system is monitored by only a few sensors,  
68       which are prone to noise, characterizing it as partially observed.

69     Unfortunately, addressing these challenges complicates model calibration.  
70     Effective model-based control necessitates accurate predictions of system be-  
71     havior, particularly ammonium concentration, over a 24-hour period, ac-  
72     counting for inherent uncertainties. This raises two critical questions:

- 73     • **Model Calibration:** how can the parameters of the model and un-  
74       certainties be determined?
- 75     • **Data Assimilation:** how can the hidden or latent space of the model  
76       be reconstructed and retrieved?

77     While existing literature often focuses on either parameter estimation or  
78     state reconstruction, there are few comprehensive studies that tackle both  
79     challenges simultaneously. Some theoretical papers explore these issues in  
80     linear cases (Holmes, 2013) or nonlinear cases (Kantas et al., 2015), but prac-  
81     tical applications are lacking. Moreover, most studies present a single model  
82     without comparing different approaches, or they focus solely on comparing  
83     inference algorithms.

84 In wastewater treatment, models like the Activated Sludge Model No.  
85 1 (ASM1) are widely used to describe nitrification processes, but parameter  
86 tuning is often neglected due to the complexity and number of parameters in-  
87 volved (Henze et al., 2006). Instead, efforts typically focus on reconstructing  
88 the latent state, such as unmeasured species concentrations, based on the  
89 few available measurements. To make calibration more manageable, some  
90 researchers simplify the physical models, reducing the number of parame-  
91 ters to be estimated (Chachuat, 2001). However, even with simplifications,  
92 accurately capturing system dynamics and uncertainties remains challeng-  
93 ing (Boulkroune, 2008). Other studies (Gernaey et al., 2004; Sánchez et al.,  
94 2002) have developed linear models from ASM1, but these models are only  
95 effective for very short-term predictions (more or less 2 hours).

96 Data-driven methods, like neural networks, learn directly from data with-  
97 out requiring explicit parameter calibration, yet they often struggle with ac-  
98 curately estimating uncertainties and adapting to changing conditions (Alvi  
99 et al., 2023, 2022). This complexity limits their integration into data assim-  
100 ilation frameworks, which are vital for estimating uncertainties from both  
101 model predictions and observations. In meteorology and oceanography, re-  
102 searchers initially relied on physical models but later incorporated data-  
103 driven techniques for better performance. However, challenges with sparse  
104 and noisy data complicate direct training (Bocquet, 2023). Recent advance-  
105 ments suggest that combining physical models with data-driven approaches  
106 can leverage the strengths of both, enabling better state reconstruction and  
107 adaptability, although issues with well-calibrated uncertainty estimates per-  
108 sist (Bonavita, 2024; Lam et al., 2023; Lguensat, 2023). In the context of

109 wastewater treatment, hybrid modeling (HM), which combines mechanistic  
110 models with data-driven approaches, has shown potential to improve pre-  
111 diction accuracy. For example, Verhaeghe et al. (Verhaeghe et al., 2024)  
112 developed a parallel hybrid model for water and resource recovery facilities  
113 (WRRFs), where a neural network was trained on the residuals of a mech-  
114 anistic model for effluent nitrate. While such hybrid approaches improve  
115 model predictions, one limitation remains: the inability to effectively quan-  
116 tify uncertainty arising from both the mechanistic and data-driven compo-  
117 nents, as well as from the observations themselves. In another direction, in  
118 meteorology and oceanography, some methods experimented using local lin-  
119 ear regression improving performance by reducing complexity (Chau et al.,  
120 2021).

121 Grey-box models combine data-driven and physics-based approaches, ef-  
122 fectively integrating both methodologies. One method involving reduced  
123 physical models to estimate dynamics is effective if simplifications are valid  
124 (Gómez-Quintero and Queinnec, 2002; Lukasse et al., 1999; Stentoft et al.,  
125 2018). In wastewater treatment plants (WWTPs), researchers like Stentoft  
126 et al. (Stentoft et al., 2018) and Lukasse et al. (Lukasse et al., 1999) have uti-  
127 lized these models; however, many studies overlook uncertainty calibration.  
128 For instance, Stare et al. (Stare et al., 2006) and Yang et al. (Yang et al.,  
129 2021) employ grey-box models to simulate the nitrification process. How-  
130 ever, their main focus is on integrating new data to estimate time-varying  
131 parameters or unobserved concentrations, rather than explicitly quantifying  
132 uncertainties for the construction of confidence intervals. As a result, they do  
133 not calibrate both model parameters and uncertainties simultaneously, but

134 instead determine uncertainty values through a posteriori tuning. Stentoft  
135 et al. (Stentoft et al., 2018) stand out for successfully calibrating both model  
136 parameters and uncertainties. Another approach, primarily explored in mete-  
137 orology and oceanography, simplifies data-driven models by integrating data-  
138 driven techniques into specific model components or throughout the data  
139 assimilation process (Fablet et al., 2021; Sacco et al., 2024), though these  
140 often rely on spatio-temporal data characteristics not directly applicable to  
141 WWTP dynamics.

142 This overview highlights a significant gap in the literature regarding the  
143 joint solution of model calibration and data assimilation. To address these  
144 gaps, this study aims to explore and compare a range of models, from purely  
145 data-driven to physics-based (mechanistic), to determine the most effective  
146 approach for 24-hour prediction of ammonium, limited to low-data-regime  
147 methods. The integration of data assimilation techniques is also pursued to  
148 estimate model parameters, uncertainties, and latent variables, thereby pro-  
149 viding a robust framework for uncertainty quantification in WWTPs. In this  
150 study, an ammonium ( $NH_4^+$ ) sensor will be continuously available, and an  
151 oxygen sensor will be provided for one experimentation. Although the focus  
152 is on WWTPs characterized by a medium-dimensional latent space, high in-  
153 teraction, and uncertainty, as well as partial observation, the methodologies  
154 developed in this study can be extended to other complex industrial systems  
155 with unknown models and uncertainties. The study first presents a detailed  
156 overview of the proposed models and the selected calibration techniques.  
157 Results are discussed, starting with simulated data to validate the models,  
158 followed by real-world data from a French wastewater treatment plant to

159 assess practical applicability.

## 160 **2. Materials and methods**

### 161 *2.1. Model's choice*

162 To compare different approaches, models are categorized into three types:  
163 white-box, grey-box, and black-box. Before detailing these categories, the  
164 general structure of the models is outlined. Assume that  $N$  data points have  
165 been collected, denoted as  $\forall i \in [1, \dots, N]$ :

- 166 •  $S[i] = (s_1[i], \dots, s_{n_x}[i]) \in \mathbb{R}^{n_x}$ , representing the state of the system at  
167 time  $t_i$ , which, in this study is the species concentrations,
- 168 •  $Y[i] = (y_1[i], \dots, y_{n_y}[i]) \in \mathbb{R}^{n_y}$ , representing the observations of the  
169 system at time  $t_i$ , which, in this study is the data collected from sensors,
- 170 •  $U[i] = (u_1[i], \dots, u_{n_u}[i]) \in \mathbb{R}^{n_u}$ , representing the control signals at time  
171  $t_i$ , which, in this study is the known aeration control,
- 172 •  $E[i] = (e_1[i], \dots, e_{n_e}[i]) \in \mathbb{R}^{n_e}$ , representing the measured exogenous  
173 signals at time  $t_i$ , which, in this study is the inflow and inlet concen-  
174 trations.

175 To separate observational noise from model noise in data assimilation, a  
176 discrete-time state space formulation is employed with two equations: the  
177 observation equation and the transition equation. By denoting  $\mathcal{N}(\mu, \Sigma)$  as  
178 the multivariate Gaussian distribution with mean  $\mu$  and covariance  $\Sigma$ , these  
179 equations can be expressed as:

$$S[i + 1] = \mathcal{M}_i(S[i], U[i], E[i]) + w_i \quad (1a)$$

$$Y[i] = H_i(S[i]) + v_i \quad (1b)$$

$$v_i \sim \mathcal{N}(0, \Sigma_\epsilon), \quad w_i \sim \mathcal{N}(0, \Sigma_m) \quad (1c)$$

180 where  $\Sigma_\epsilon$  and  $\Sigma_m$  are covariance matrices for the observation noise and the  
 181 model noise, respectively. The timestep of the model and observation can  
 182 differ, but for the sake of readability, it is considered to be the same. This  
 183 state space approach allows for modeling any dynamical system generically,  
 184 representing any system in this form (Durbin and Koopman, 2012). The  
 185 purpose of this section is to define the functions  $\mathcal{M}_i$  and  $H_i$  for each model.  
 186 For the rest of the paper,  $\Sigma_\epsilon$  and  $\Sigma_m$  are considered as diagonal matrices.  
 187 However, the methodology can be extended to alternative covariance struc-  
 188 tures, including time-dependent formulations influenced by various factors  
 189 (Cheng et al., 2023).

### 190 2.1.1. White-box model

191 The first model, introduced by Chachuat et al. (Chachuat, 2001) and  
 192 used in various studies (Boukroune, 2008; Chachuat, 2001), is a simplified  
 193 version of the ASM1 model (Henze et al., 2006). This model, specifically  
 194 developed for AASP, removes slow dynamics such as sludge weight variation  
 195 and comprises five states:

$$S[i] = (x_{COD}[i], s_O[i], s_{NO}[i], s_{NH}[i], s_{ND}[i]).$$

196 These states represent the mass of COD (Chemical Oxygen Demand) and  
 197 the concentrations of oxygen ( $O$ ), nitrite/nitrate ( $NO$ ), ammonium ( $NH$ ),

198 and organic nitrogen ( $ND$ ), respectively. The model maintains the structure  
 199 of ASM1, which describes nitrogen processing in a Continuously Stirred Tank  
 200 Reactor (CSTR) with a constant volume  $V$  and a continuous inflow rate  
 201  $Q_{in}[i]$ . Using the notation  $\dot{S}[i]$  to define the time derivative of  $S[i]$ , the  
 202 concentration evolution within the reactor is described by:

$$\dot{S}[i] = \frac{Q_{in}^{in}[i]}{V}(S^{in}[i] - S[i]) + r(i, S[i], U[i]) + w_i \quad (2)$$

203 where  $E[i] = (Q_{in}[i], S^{in}[i])$  denotes the inflow and inlet water concentra-  
 204 tions with  $S^{in}[i] = (s_{NH}^{in}[i], \dots, s_{ND}^{in}[i])$ , and  $U[i] \in \{0, 1\}$  represents the  
 205 aeration control. The ASM1 model mathematically represents the function  
 206  $(i, S, U) \mapsto r(i, S[i], U[i])$ , detailing the reactions in the reactor tank. De-  
 207 tailed descriptions of this function can be found in the references (Henze et al.,  
 208 2006, 2008; Metcalf & Eddy Inc. et al., 2013). Chachuat et al. (Chachuat,  
 209 2001) simplified the latent space  $S[i]$  from 13 to 5 states and modified the  
 210 structure of the function  $r(i, S[i], U[i])$ .

211 For the remainder of this article, this model will be referred to as the  
 212 white-box model. Our objective is to estimate the model’s parameters, ini-  
 213 tial conditions, and the concentrations within the latent space  $S(t)$  to closely  
 214 replicate the dynamics observed in the actual tank. However, as detailed in  
 215 the previous section, this is not a straightforward task. In current literature,  
 216 no attempts have been made to estimate all these components simultane-  
 217 ously using this model. For example, Boulkroune et al. (Boulkroune, 2008)  
 218 and Chachuat et al. (Chachuat, 2001) estimated only a subset of model’s  
 219 parameters, fixing others within the latent space to simplify the problem.  
 220 Moreover, as identified by Boulkroune et al. (Boulkroune, 2008), the model

221 is only observable under certain conditions, particularly when variables such  
222 as oxygen, nitrate, and ammonium concentrations are being monitored.

### 223 2.1.2. Grey-box model

224 To further simplify the model and enhance observability, some research  
225 has focused specifically on the dynamics of interest (Bertret et al., 2024;  
226 Lukasse et al., 1999; Stentoft et al., 2018). For instance, if the primary  
227 concern is the concentration of ammonium ( $s_{NH}$ ), the latent space can be  
228 reduced to include only ammonium concentration. This approach targets  
229 the relevant dynamics and eliminates dependencies on unmeasured variables.  
230 The models proposed by Lukasse et al. and Stentoft et al. (Lukasse et al.,  
231 1999; Stentoft et al., 2018) are based on ammonium and nitrate concentra-  
232 tions, while recent work focused solely on ammonium (Bertret et al., 2024).  
233 In our case, given the primary focus on ammonium, the model considered in  
234 this section will be the one proposed in recent work (Bertret et al., 2024),  
235 defined as:

$$\dot{s}_{NH}[i] = \frac{Q^{in}[i]}{V} [s_{NH}^{in}[i] - s_{NH}[i]] - \beta U[i] \frac{s_{NH}[i]}{s_{NH}[i] + K} + w_i \quad (3)$$

236 where  $K$ ,  $\beta$ , and  $V$  are model's parameters. This model will be referred to as  
237 the grey-box model as it's a very simplified physical model where the model's  
238 parameters have no more physical meaning and will be determined by data.

239 However, while most studies attempt to estimate all the model's param-  
240 eters (Lukasse et al., 1999; Stentoft et al., 2018), few consider calibrating  
241 uncertainties. In the literature, to the best of our knowledge, only Lukasse

242 et al. (Lukasse et al., 1999) and Bertret et al. (Bertret et al., 2024) have  
243 addressed uncertainty calibration.

### 244 2.1.3. *Black-box model*

245 Finally, the first two approaches are compared to models that are purely  
246 data-driven. Two types of data-driven models are explored. The first one  
247 involves using a linear state space model where the latent space dimension is  
248 equal to the number of observed variables:

$$S[i + 1] = S[i] + AS[i] + BU[i] + CE[i] + w_i. \quad (4)$$

249 The objective of this model is to determine whether the dynamics of the  
250 system can be adequately represented by a linear model, as suggested in the  
251 literature (Gernaey et al., 2004; Sánchez et al., 2002). It will be called linear  
252 black-box model.

253 The second data-driven model is inspired by the recent work of Chau  
254 et al. (Chau et al., 2021; Lguensat et al., 2017). In this case, a local linear  
255 regression (LLR) is employed to model the reactor’s dynamics. This approach  
256 has demonstrated strong performance when the latent space is small (Chau  
257 et al., 2021), effectively reconstructing the hidden states. Mathematically,  
258 this model is expressed as:

$$S[i + 1] = S[i] + \text{LLR}(S[i], U[i], E[i]) + w_i \quad (5)$$

259 where LLR represents the output of the local linear regression. Detailed  
260 information on the methods and hyperparameters can be found in Chau et  
261 al. (Chau et al., 2021; Lguensat et al., 2017). Not going into the details, the

262 methods identify the  $k$  nearest neighbors of the current point, make a linear  
263 regression on the selected point and predict the projection of the current  
264 point using the linear regression. It will be called non-parametric black-box  
265 model. The advantages of these two black-box methods are their proven  
266 utility in parameter fitting and latent space reconstruction. However, the  
267 second method has not yet been applied to wastewater treatment plants and  
268 may offer a promising middle ground between linear models and artificial  
269 neural networks (ANN) when data availability is limited. Other models,  
270 such as artificial neural networks or Gaussian processes (GP), could also be  
271 considered within the black-box approach. However, LLR can be viewed as a  
272 sparse GP (Stulp and Sigaud, 2015), and in scenarios with limited data, ANN  
273 tends to generalize less effectively than LLR (see Supplementary Material and  
274 (Otchere et al., 2021)).

275 This section has presented the different types of models. However, the  
276 simultaneous estimation of the model's parameters and reconstruction of the  
277 latent space has not yet been addressed. This will be the focus of the next  
278 section.

## 279 *2.2. Model's calibration*

280 In the context of data assimilation, several challenges arise concerning  
281 inference (Durbin and Koopman, 2012):

- 282 • reconstructing hidden states (the focus of most studies);
- 283 • calibrating model's parameters and uncertainties (addressed in fewer  
284 studies).

285        There are various solutions to the first challenge. Variational methods and  
286 sequential methods are commonly employed (Carrassi et al., 2018). Sequen-  
287 tial methods, including the Kalman Filter, were developed first. However,  
288 these methods can become computationally expensive for high-dimensional  
289 systems, which led to the development of variational methods (Sasaki, 1970).  
290 Variational methods are primarily used in high-dimensional systems, such as  
291 those found in oceanography or meteorology, and work by minimizing a cost  
292 function (Carrassi et al., 2018; Sasaki, 1970). This study focuses on systems  
293 of medium to small dimensions, emphasizing sequential methods, which offer  
294 stronger theoretical guarantees.

295        For linear systems, the Kalman Filter and Kalman Smoother provide a  
296 closed-form solution. However, in nonlinear systems, adaptations are neces-  
297 sary. These adaptations fall into two categories:

- 298        • Linearization techniques, such as the Extended Kalman Filter (EKF)  
299        (Anderson and Moore, 2012) and the Unscented Kalman Filter (UKF).
- 300        • Sequential Monte Carlo methods, including the Ensemble Kalman Fil-  
301        ter (EnKF) (Anderson and Moore, 2012; Evensen and Leeuwen, 2000)  
302        and the Particle Filter (PF) (Doucet and Johansen, 2011; Kantas et al.,  
303        2015).

304        While linearization methods can be faster, they may not be suitable when  
305 non-linearities between timesteps are significant, leading to less accurate so-  
306 lutions, as indicated by various studies (Dreano et al., 2017; Stentoft et al.,  
307 2018). Sequential Monte Carlo methods (Doucet et al., 2001; Doucet and  
308 Johansen, 2011) estimate the target distribution using particles. The EnKF

309 and Ensemble Kalman Smoother (EnKS) assume that the state is Gaussian  
 310 at each timestep, while Particle Filters and their associated smoothers are  
 311 not restricted to Gaussian distribution. Although Particle Filters is based on  
 312 theoretical results that provide more generality and flexibility than the EnKF,  
 313 they usually require more particles than the EnKF for equivalent performance  
 314 in estimating the posterior mean, making them slower (Chau et al., 2021).  
 315 To address this, methods like the Conditional Particle Filter (CPF) have  
 316 been developed to retain the benefits of Particle Filters while reducing the  
 317 number of particles needed (Lindsten, 2013; Lindsten et al., 2014). The CPF  
 318 involves constructing a Markov kernel whose asymptotic law is the posterior  
 319 distribution of the hidden state. Instead of increasing the number of parti-  
 320 cles, this method iterates the kernel to obtain an accurate approximation of  
 321 the distribution.

322 In summary, various tools for addressing the problem of reconstructing  
 323 hidden states have been explored. However, the calibration of model paramete-  
 324 rs and uncertainties remains to be addressed. The goal is to calibrate these  
 325 parameters in order to maximize the likelihood of the observations, expressed  
 326 as:

$$L(Y; \theta) = p(Y[1], \dots, Y[n]; \theta) = p(Y[1]) \prod_{i=2}^N p(Y[i] | Y[1 : i - 1]; \theta) \quad (6)$$

327 where  $Y[1 : i - 1] = (Y[1], \dots, Y[i - 1])$ . Almost any problem can be  
 328 framed as a differentiable program. Thus, one potential approach is to differ-  
 329 entiate Sequential Monte Carlo methods to find the model's parameters that  
 330 maximize the likelihood (Durbin and Koopman, 2012; Stentoft et al., 2018).

331 This approach, however, can be unstable and may result in challenging op-  
332 timization problems. An alternative is to maximize the complete likelihood,  
333 incorporating both observations and hidden states. Since the hidden states  
334 are unknown, a common solution is to use the 'Expectation-Maximization'  
335 (EM) algorithm. This algorithm divides the problem into two steps: the 'E-  
336 Step,' where the posterior distribution of the hidden state is reconstructed  
337 using filters and smoothers, and the 'M-Step,' where the complete likelihood  
338 is optimized using observations and the predicted hidden states from the  
339 E-Step (Chau et al., 2021; Dreano et al., 2017; Lindsten, 2013). This ap-  
340 proach is known as the Expectation-Maximization (EM) algorithm for linear  
341 systems (Dempster et al., 1977; Durbin and Koopman, 2012). Extensions of  
342 this method using Sequential Monte Carlo methods for the "E-Step", such  
343 as Stochastic Approximation Expectation-Maximization (SAEM) and Monte  
344 Carlo Expectation-Maximization (MCEM), have been developed for nonlin-  
345 ear systems (Celeux et al., 1995; Lindsten, 2013).

346 Chau et al. (Chau et al., 2021) compared EKF, EnKF, PF, CPF, and  
347 associated smoothers for parameter estimation in nonlinear models. Based  
348 on their conclusions, the focus will be on the Particle Filter (PF) and Con-  
349 ditional Particle Filter (CPF) associated with the Backward Smoother (BS)  
350 for nonlinear models due to their efficiency, generality, and performance. For  
351 linear models, the Kalman Filter and Smoother will be employed. These  
352 smoothers and filters will be combined with the EM or MCEM methods,  
353 depending on whether the model is linear or nonlinear, as these methods  
354 have demonstrated strong performance in parameter calibration (Chau et al.,  
355 2021). Table 1 summarizes the methods chosen for each model presented in

356 the previous section.

Table 1: Summary of calibration methods for different model types in data assimilation. The table presents the approaches used for hidden state estimation and parameter calibration across various models considered in this study. The number of parameters for each model, including both model parameters and uncertainties, is also provided, where  $n_Y$  represents the number of observations.

| Model                    | Hidden State Estimation    | Parameter Calibration | # parameters    |
|--------------------------|----------------------------|-----------------------|-----------------|
| Linear Black-Box         | Kalman Filter and Smoother | EM                    | $2n_Y^2 + 4n_Y$ |
| Grey-Box                 | PF and BS                  | MCEM                  | 5               |
| White-Box                | CPF and BS                 | MCEM                  | $24 + n_Y$      |
| Non-Parametric Black-Box | CPF and BS                 | MCEM                  | $2n_Y^2 + 4n_Y$ |

357

358 The methods outlined in this section allow us to reconstruct hidden states  
359 and calibrate model’s parameters and uncertainties, a task that is rarely un-  
360 dertaken in the literature on wastewater treatment plants. These methods  
361 enable the calibration of any type of state-space model, followed by the refine-  
362 ment of hidden state estimates using observations. This approach is more  
363 robust and powerful than standard Prediction Error Minimization (PEM)  
364 methods (Ljung, 1998), as demonstrated in recent work (Bertret et al., 2024).

### 365 3. Results and discussion

366 As detailed, the primary objective of this study is to develop a model  
367 capable of making rolling predictions of ammonium concentration. Various  
368 strategies, ranging from purely physical models to data-driven models, have  
369 been presented. To achieve this, certain physical aspects of the system are

370 simplified to isolate the most relevant dynamics for modeling. These simpli-  
371 fications are necessary to ensure tractability and solvability given the limited  
372 data available, without compromising the system’s essential chemical prop-  
373 erties. Consequently, the models are trained on sparse datasets, such as  
374 one week of observations. This approach is driven by two key factors: first,  
375 simplifying the models by assuming constant concentrations over short time  
376 periods, such as the size of the dataset, allows for calibration on recent data;  
377 second, recent data more accurately captures the system’s immediate dy-  
378 namics, particularly when external conditions, such as influent composition  
379 or temperature, change rapidly. This strategy reflects practical constraints,  
380 where extensive experimentation and long-term data collection are often in-  
381 feasible, highlighting the need for adaptable models that perform well with  
382 minimal data.

### 383 *3.1. Simulation results*

384 This section aims to compare all the presented models and demonstrate  
385 how data assimilation can effectively handle latent space estimation and pro-  
386 vide uncertainty quantification. First, the models are compared on simulated  
387 data, and then the best-performing model is tested in a real installation  
388 in France in section 3.2. The analysis was conducted using the Julia pro-  
389 gramming language, and the experimental code is available <sup>1</sup>. As previously  
390 mentioned, real wastewater treatment plants are typically monitored with a  
391 limited number of noisy sensors. Most municipal plants are equipped with

---

<sup>1</sup>The code for reproducing results of this paper can be accessed at <https://github.com/vbertret/WWTP-DataAssimilation-ModelComparison>.

392 only one oxygen and/or redox sensor, and occasionally an ammonium sensor,  
393 all placed at a single location in the reactor tank. Consequently, the sensor  
394 readings may be noisy and might not accurately represent the conditions  
395 throughout the entire reactor. This study focuses on quantifying ammonium  
396 concentration, assuming the availability of an ammonium sensor, along with  
397 an additional oxygen sensor in some cases.

### 398 3.1.1. *Experimental Setup*

399 The primary goal of this study is to compare the models in terms of  
400 24-hour horizon prediction, which is crucial for effective control. A fair com-  
401 parison is ensured by testing the models under various conditions, including  
402 different types of training and test data, to provide robust results. To achieve  
403 this, the ASM1 model, as described earlier, is used to generate synthetic  
404 data. As previously noted, the estimation of inlet concentrations is left for  
405 a future study, so inlet concentrations and inflow data from the Benchmark  
406 Simulation Model (BSM) are utilized. The control of the aeration system  
407 is determined by a strategy resembling ORP (Oxydo Reduction Potential)  
408 control, where aeration is activated or deactivated based on thresholds for  
409  $NH_4^+$  and  $NO_3^-$  (Paul et al., 1998). To make the simulation results more  
410 realistic, noise is added to the observations, and training is performed on 50  
411 independent realizations of the observation time series. The system is first  
412 simulated for 20 days to reach a steady state, followed by a training period  
413 ( $N_{training}$  points) and a rolling test period ( $N_{test}$  points), with observations  
414 collected every 5 minutes.

415 The test period lasts 14 days, during which 24-hour horizon predictions  
416 are made iteratively. After each 1-hour interval, new observations are assim-

417 ilated, and the predictions are updated for the next 24 hours. This approach  
 418 allows us to evaluate the models across a wide range of inflow and concen-  
 tration conditions as illustrated on the Figure 2.

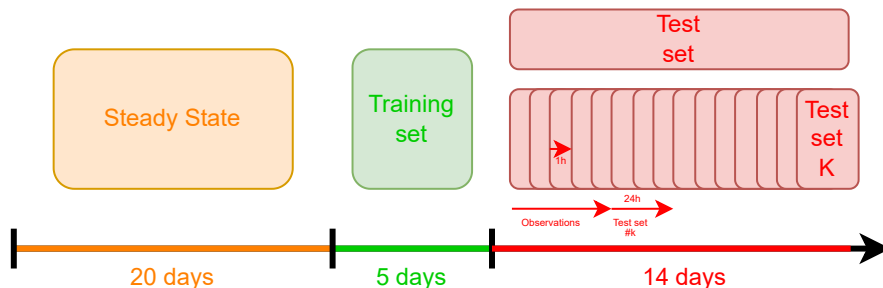


Figure 2: Dataset construction used in model training and testing. The system is first simulated to reach a steady state, followed by a training period of 5 days and a rolling test period of 14 days. Observations are assimilated every hour, and 24-hour horizon predictions are made iteratively.

419

### 420 3.1.2. Metrics

421 In this study, three distinct metrics are used to evaluate the model's per-  
 422 formance: Root Mean Squared Error (RMSE), Average Width (AW), and  
 423 Coverage Probability (CP). Each metric serves a specific purpose in assess-  
 424 ing different aspects of the model's predictions, providing a comprehensive  
 425 evaluation of both accuracy and reliability. Let  $\hat{Y}[i|i-h]$ ,  $\hat{Y}_{LB}^\alpha[i|i-h]$ , and  
 426  $\hat{Y}_{UB}^\alpha[i|i-h]$  represent, respectively, the mean prediction, the lower bound,  
 427 and the upper bound of the model's  $\alpha\%$  confidence interval at time  $i$ , given  
 428 that the system is observed up to time  $i-h$ . Since predictions can be made at  
 429 different horizons, all metrics are computed over a range defined by a starting  
 430 horizon  $H_1$  and an ending horizon  $H_2$  (with  $H_2 > H_1$ ), and the results are

431 averaged over this range.

432 *RMSE (Root Mean Squared Error)*. This metric evaluates the accuracy of  
 433 the model by measuring the difference between the predicted and true values.  
 434 Specifically, it quantifies how well the model predicts concentrations at each  
 435 time step. A lower RMSE indicates that the model is closer to the actual  
 436 system behavior. For rolling predictions, the RMSE is computed as:

$$\mathcal{L}_{RMSE}(H_1, H_2) = \frac{1}{H_2 - H_1} \sum_{h=H_1}^{H_2} \sqrt{\sum_{i=H_2}^{N_{test}} \frac{\|\hat{Y}[i|i-h] - Y[i]\|_2^2}{N_{test} - H_2}} \quad (7)$$

437 RMSE provides a direct measure of the model’s predictive accuracy by  
 438 focusing on the deviation between predicted and observed concentrations.

439 *AW (Average Width)*. This metric evaluates the precision of the model’s  
 440 predictions by assessing the width of the confidence intervals. A narrower  
 441 interval indicates that the model is more certain about its predictions, while  
 442 a wider interval implies more uncertainty. For rolling predictions, the AW is  
 443 computed as:

$$\mathcal{L}_{AW}(H_1, H_2) = \frac{1}{H_2 - H_1} \sum_{h=H_1}^{H_2} \sum_{i=H_2}^{N_{test}} \frac{[\hat{Y}_{UB}^\alpha[i|i-h] - \hat{Y}_{LB}^\alpha[i|i-h]]}{N_{test} - H_2} \quad (8)$$

444 AW assesses the model’s confidence in its predictions. A smaller interval  
 445 indicates higher precision, while still accounting for uncertainty.

446 *CP (Coverage Probability)*. The CP metric evaluates the reliability of the  
 447 confidence intervals by checking how frequently the true concentration values  
 448 fall within the predicted intervals. Ideally, the actual values should fall within

449 the interval as frequently as the interval was calibrated to (e.g., 95%). For  
 450 rolling predictions, the CP is computed as:

$$\mathcal{L}_{CP}(H_1, H_2) = \frac{1}{H_2 - H_1} \sum_{h=H_1}^{H_2} \sum_{i=H_2}^{N_{test}} \frac{\mathbf{1}_{\{\hat{Y}_{LB}^\alpha[i|i-h] \leq Y[i] \leq \hat{Y}_{UB}^\alpha[i|i-h]\}}}{N_{test} - H_2} \quad (9)$$

451 CP measures how well-calibrated the model is in terms of its uncertainty  
 452 estimates, ensuring the intervals contain the correct proportion of true values.

### 453 3.1.3. Forecasting performance

454 The primary goal of this experiment is to evaluate the forecasting per-  
 455 formance of each model when only ammonium concentration is observed, re-  
 456 flecting a real-world scenario where sensor availability is limited due to costs  
 457 and calibration requirements. The models were trained for 5 days using a  
 458 fixed inlet concentrations and then tested over a 14-day period, as described  
 459 earlier. Figure 3 illustrates the performance metrics across various prediction  
 460 horizons.

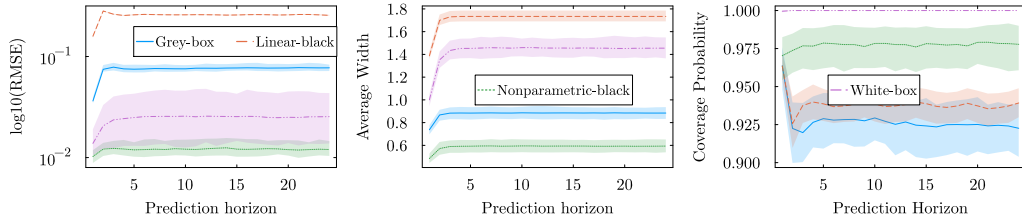


Figure 3: Comparison of model performance metrics. This figure presents the root mean squared error (RMSE), average width (AW), and coverage probability (CP) for different models across various prediction horizons (in hours). The value at  $x = 5h$  corresponds to the average metric computed over the interval  $\mathcal{L}_{\text{RMSE}}(4h, 5h)$ , and similarly for the other metrics. The solid lines represent the mean values, while the shaded areas denote the 95% confidence intervals.

461 The results indicate that the linear model performs the worst in terms of  
462 RMSE, corroborating findings from previous studies (Gernaey et al., 2004;  
463 Sánchez et al., 2002). In contrast, the non-parametric black-box model sig-  
464 nificantly outperforms all others across all metrics. This counter-intuitive  
465 result will be discussed later, as it may seem surprising given that grey-box  
466 and white-box models contain more information about the system. The com-  
467 parison between the white-box and grey-box models is less straightforward:  
468 while the white-box model shows a lower RMSE, its confidence intervals are  
469 considerably wider, suggesting higher variability and instability. This insta-  
470 bility likely stems from the white-box model being initialized with parameter  
471 values close to those used in the simulator, resulting in good RMSE but poor  
472 training stability. Interestingly, no error accumulation is observed under  
473 these controlled conditions, likely due to the use of a simulator with limited  
474 and controlled sources of uncertainty. However, in real-world applications,  
475 error accumulation may arise due to higher uncertainty relative to ammonium  
476 dynamics as shown in Section 3.2. As seen in Figure 3, all models are well-  
477 calibrated with a coverage probability close to 95%, except for the white-box  
478 model, which suffers from overly broad confidence intervals. This finding  
479 underscores the effectiveness of our parameter and uncertainty calibration  
480 approach. Table 2 provides detailed information on the estimated standard  
481 deviations for the models and observations, as well as training times.

482 The linear model is the fastest to train, followed by the grey-box and non-  
483 parametric models. Although the non-parametric model requires 20 times  
484 more training time than the linear model and 5 times more than the grey-  
485 box, the difference remains manageable. However, the white-box model is

Table 2: Model training time in seconds and standard deviation estimates in mg/L. This table summarizes the training time, estimated standard deviation of observations, and estimated standard deviation of the model for each of the four models. The known standard deviation of observations was set to **0.2**.

| Model               | Training Time<br>(seconds) | Observation Standard<br>Deviation (mg/L) | $s_{NH}$ Standard<br>Deviation (mg/L) |
|---------------------|----------------------------|--|---------------------------------------|
| Grey-box            | 8.13e+02                   | 1.84e-01                                 | 1.03e-01                              |
| Linear-black        | <b>1.62e+02</b>            | 1.62e-01                                 | 1.80e-01                              |
| Nonparametric-black | 3.99e+03                   | <b>1.97e-01</b>                          | <b>6.23e-02</b>                       |
| White-box           | 2.94e+05                   | 1.89e-01                                 | 1.09e-01                              |

486 extremely time-consuming, taking 1850 times longer than the linear model  
 487 and 375 times longer than the grey-box, making it impractical for real-world  
 488 use in parameter and uncertainty calibration.

489 Regarding the estimation of observation noise, all models perform well,  
 490 with values close to the true standard deviation of 0.2. Notably, the non-  
 491 parametric model achieves a value of 0.197 and also shows the lowest standard  
 492 deviation for the model itself. As highlighted in Table 2, the non-parametric  
 493 model outperforms the others overall. Figure 4 showcases a single 24-hour  
 494 prediction realization for each model.

495 The linear-black model, as expected, struggles with complex dynamics,  
 496 such as saturation near zero concentration. The grey-box model delivers ac-  
 497 curate predictions and reasonable confidence intervals, though it tends to be  
 498 biased at high ammonium levels, possibly due to oversimplifications in the  
 499 physical model. This bias is somewhat compensated by wider confidence  
 500 intervals. The white-box model is well-calibrated in terms of mean predic-

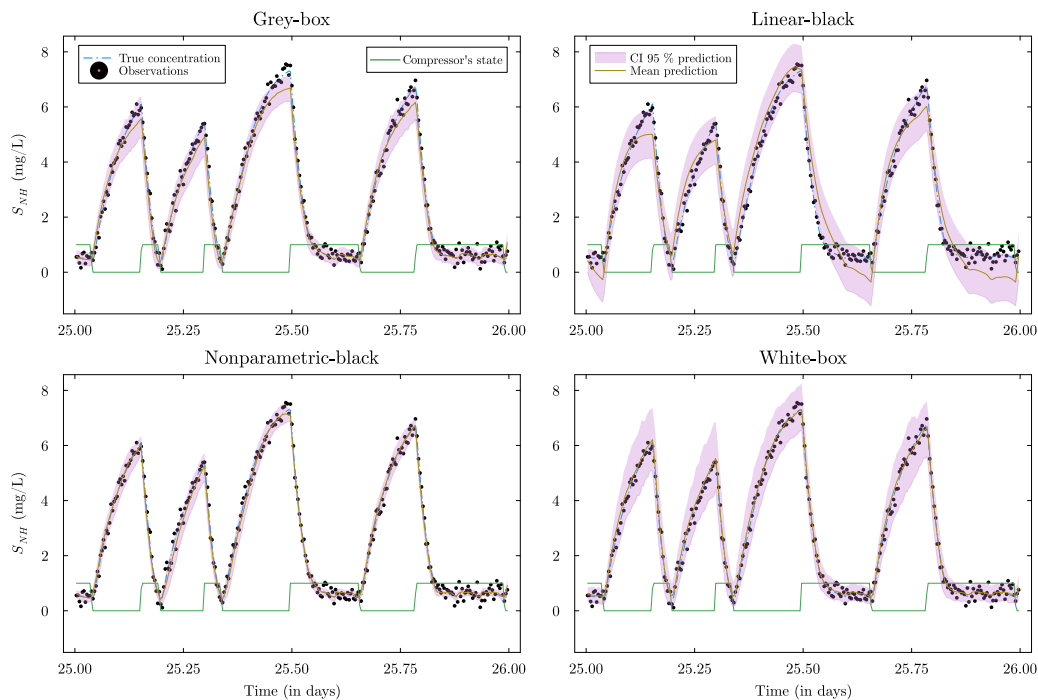


Figure 4: 24-hour prediction realizations for each model. This figure compares the predicted ammonium concentrations (solid lines) with the true concentrations (dots) for one realization of each model. The shaded areas represent the 95% confidence intervals.

501 tions but suffers from excessively broad confidence intervals, reflecting its  
 502 instability. Lastly, the non-parametric black-box model demonstrates strong  
 503 performance in both mean predictions and confidence intervals, effectively  
 504 capturing the information from the training set.

505 In summary, the linear black-box model is fast and well-calibrated in  
 506 terms of uncertainty but lacks the flexibility needed to capture the complex  
 507 dynamics of a wastewater treatment reactor. The white-box model, while  
 508 potentially powerful, is difficult to train, time-consuming, and highly depen-  
 509 dent on good initial model's parameters, making it less generalizable. The

510 grey-box model offers a good balance between accuracy and computational  
 511 efficiency, though some simplifications may degrade its performance. Fi-  
 512 nally, the non-parametric black-box model excels in both predictive accuracy  
 513 and uncertainty estimation, highlighting its potential for robust modeling in  
 514 wastewater treatment applications.

### 515 3.1.4. Added value of other concentration

516 Having compared the models using only ammonium concentration data,  
 517 this section explores whether incorporating additional sensors, such as an  
 518 oxygen sensor, can enhance model performance. Based on the findings from  
 519 the previous section, the non-parametric black-box model emerged as the  
 520 best-performing model, so this analysis will focus on that model alone. The  
 521 experimental setup remains consistent with the first experiment, but now  
 522 the model has access to both oxygen and ammonium concentration data.  
 523 Table 3 presents a comparison of the model’s performance when observing  
 524 ammonium alone versus both ammonium and oxygen.

Table 3: Comparison of model performance with different observed species. This table compares the root mean squared error (RMSE), average width (AW), and coverage probability (CP) for the non-parametric black-box model when observing only ammonium ( $NH_4$ ) versus observing both ammonium and oxygen ( $NH_4, O_2$ ). The results are shown for ammonium prediction.

| Observed species | Observed Standard Deviation | $s_{NH}$ Standard Deviation | $\mathcal{L}_{RMSE}(1H, 24H)$ | $\mathcal{L}_{AW}(1H, 24H)$ | $\mathcal{L}_{CP}(1H, 24H)$ |
|------------------|-----------------------------|-----------------------------|-------------------------------|-----------------------------|-----------------------------|
| $[Y_{NH}]$       | 1.97e-01                    | <b>6.23e-02</b>             | <b>1.21e-02</b>               | <b>5.88e-01</b>             | 9.77e-01                    |
| $[Y_{NH}, Y_O]$  | 1.98e-01                    | 9.05e-02                    | 3.47e-02                      | 7.94e-01                    | 9.69e-01                    |

525 The results indicate that the observation noise standard deviation is ac-  
 526 curately estimated in both scenarios, validating the robustness of our un-

527 certainty estimation methodology. However, the model that observes only  
 528 ammonium consistently outperforms the model that observes both ammo-  
 529 nium and oxygen across all metrics. Specifically, the root mean squared error  
 530 (RMSE), confidence interval width (AW), the coverage probability (CP) and  
 531 model noise standard deviation are all lower when the model focuses solely  
 532 on ammonium.

533 This counterintuitive outcome might be explained by the additional com-  
 534 plexity introduced when the model must simultaneously reconstruct both  
 535 species. The added noise in estimating oxygen appears to interfere with the  
 536 accuracy of ammonium predictions. Figure 5 illustrates a 24-hour prediction  
 537 realization for the model observing both oxygen and ammonium.

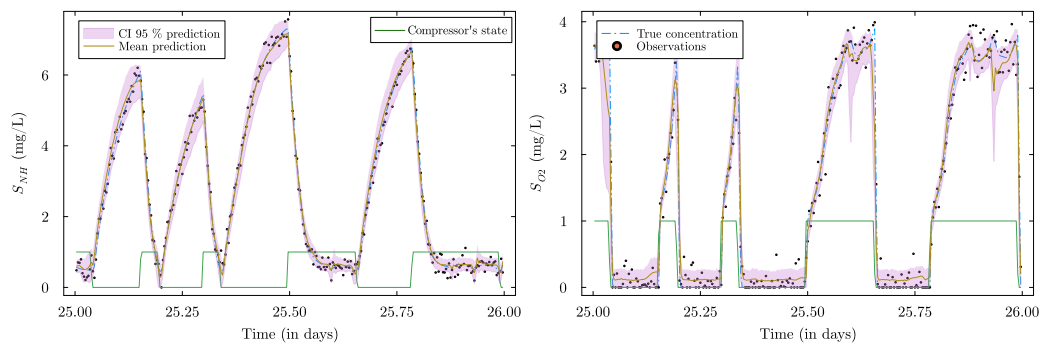


Figure 5: 24-hour prediction realizations for ammonium and oxygen. This figure compares the predicted and true concentrations of ammonium (left panel) and oxygen (right panel) for one realization of the non-parametric black-box model. The solid lines represent the mean predictions, the dotted line the true concentrations, and the shaded areas indicate the 95% confidence intervals.

538 The model successfully captures the dynamics of ammonium, though  
 539 some errors that were not present in the single-species scenario become ev-  
 540 ident. The oxygen dynamics are also well-reconstructed; however, the as-

541 sumption of Gaussian noise for oxygen may not be optimal, suggesting that  
542 exploring alternative distributions could be beneficial in future studies. In  
543 conclusion, adding an additional sensor, such as oxygen, does not appear to  
544 improve model performance when the target species (ammonium) is already  
545 being observed. This finding suggests that, at least for the non-parametric  
546 black-box model, investing in additional sensors may not be justified if the  
547 primary species of interest is already monitored. However, it would be valu-  
548 able to investigate whether observing other species could be beneficial when  
549 the target species is not directly observed—though such an exploration is  
550 beyond the scope of this paper.

### 551 *3.2. Real system experiments*

552 The non-parametric black-box model showed promising results on simu-  
553 lated data. However, to assess its performance in a real-world setting and  
554 ensure that the simulated data were not overly simplified, the model needed  
555 to be tested on real data. To achieve this, the non-parametric black-box  
556 model was applied to a real-world wastewater treatment plant located in  
557 Acigné, France, which serves a population equivalent (PE) of 14,000. Data  
558 collection spanned from April 12 to May 8, 2023, with the first five days used  
559 for training the model and the subsequent days for performance evaluation.  
560 Inflow and ammonium concentrations ( $NH_4^+$ ) were recorded every minute.  
561 Since the inlet concentrations were not directly measured, it was excluded  
562 from the exogenous signals used as model inputs.

563 Figure 6 displays four 24-hour model predictions for the real dataset, each  
564 beginning at a different time of day. Although the real concentration data is  
565 unavailable, sensor readings were used for analysis. Compared to the simula-

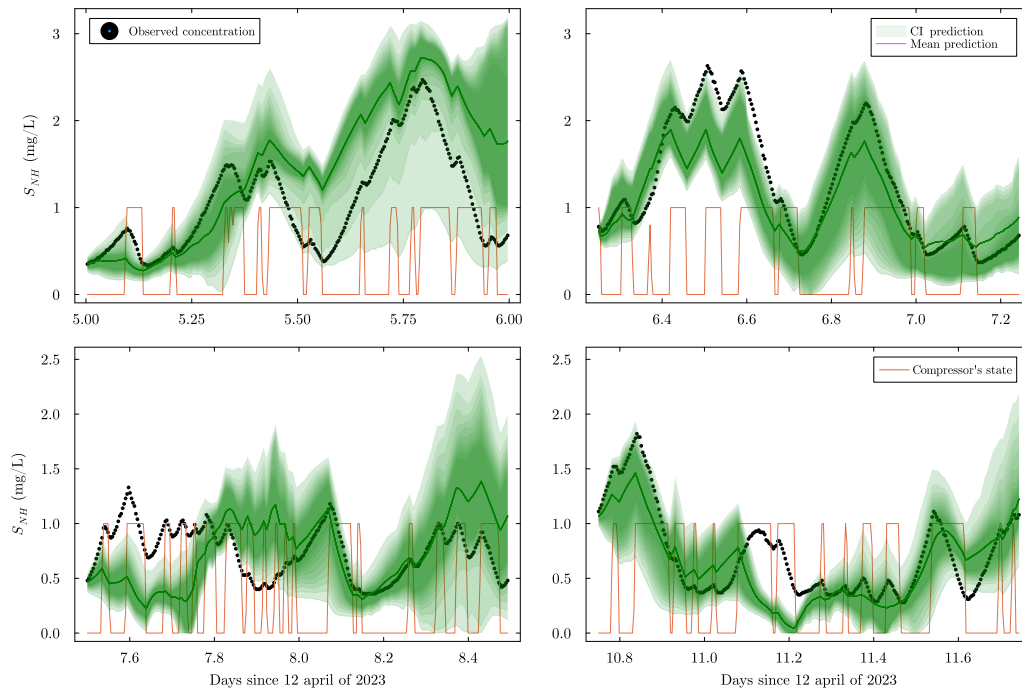


Figure 6: 24-hour model predictions on real system data, initiated at different times of the day. The black dots represent sensor readings of ammonium concentration ( $s_{NH}$ ) in mg/L, while the green-shaded regions represent different percentages of confidence interval (CI) for the predicted values. Darker shading corresponds to lower percentages of values contained within the interval.

566 tion results, the real-world dynamics are less accurately captured. The model  
 567 shows higher uncertainty, as reflected by the wider confidence intervals, which  
 568 still manage to encapsulate the majority of the observations. Notably, the  
 569 model's deviations are more pronounced in amplitude rather than in trend  
 570 evolution, likely due to the assumption—borrowed from the simulation—that  
 571 the inlet concentrations remain constant. This assumption does not hold  
 572 true for the real system, where significant daily variations in inlet concen-  
 573 trations are observed (Gerney et al., 2014). As the model lacks this critical

574 information, it must generalize across varying conditions, leading to higher  
575 confidence intervals despite generally capturing the underlying dynamics. To  
576 provide further insights into the predicted distribution, confidence intervals  
577 of varying percentages were included, highlighting the asymmetry in the in-  
578 tervals, as the distribution is non-Gaussian. For example, in the top left plot,  
579 although the 95% confidence interval is wide near the end of the prediction,  
580 the 90% interval is significantly narrower, demonstrating a key strength of  
581 the method used in this study. This behavior indicates that the model can  
582 still maintain tight predictions even when broader uncertainty exists for the  
583 outer confidence bands. In the bottom prediction, however, certain sections  
584 appear too confident, with overly narrow confidence intervals. This might  
585 indicate that the model is underestimating uncertainty during these periods.  
586 Nonetheless, it is a localized issue and doesn't detract from the overall per-  
587 formance. For instance, in the bottom right prediction, between days 11.0  
588 and 11.2, the model predicts a decrease in ammonium concentration dur-  
589 ing aeration, whereas the observed values show an increase, likely caused by  
590 an unusual spike in inlet pollution (ammonium). As this is an infrequent,  
591 abnormal event, it is understandable that the model failed to capture it.

592 To further explore the impact of inlet variation, Figure 7 compares the  
593 scaled mean derivative of the non-aeration periods across different times of  
594 the day. Figure 7 shows that under simulated conditions with a fixed inlet,  
595 the derivative remains nearly constant, as expected. However, in the real  
596 data, the derivative fluctuates throughout the day, increasing during day-  
597 light hours and decreasing at night. This diurnal variation mirrors patterns  
598 observed in other studies measuring real wastewater treatment plant opera-

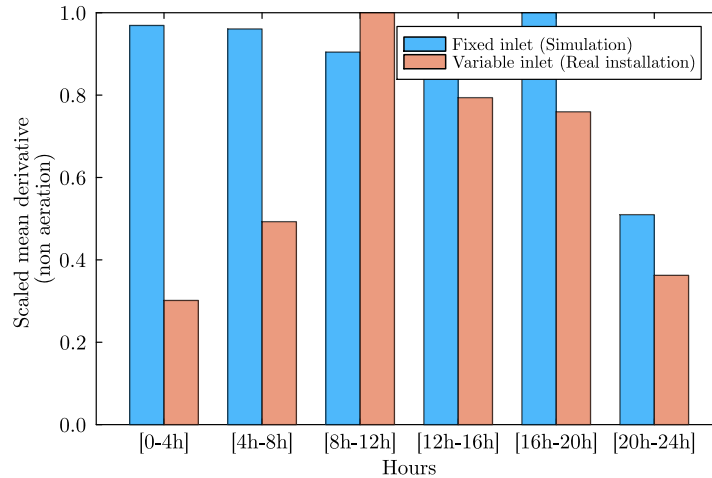


Figure 7: Scaled mean derivative of non-aeration periods comparing simulated and real installation data across different time windows. The blue bars represent fixed inlet simulations, showing nearly constant derivative values. The orange bars represent real installation data, exhibiting significant diurnal variation, with increased derivatives during the day and decreased values at night, reflecting the inlet variability observed in real wastewater treatment plants.

599 tions. These findings suggest that the model’s diminished performance in the  
 600 real case may stem from the unaccounted variability in inlet concentrations.  
 601 To improve accuracy, additional covariates such as time of day, rainfall, and  
 602 other environmental factors should be incorporated into the model. These  
 603 improvements will be addressed in future research.

#### 604 4. Conclusion

605 This study compared different types of models, from white-box to black-  
 606 box models, to predict ammonium concentration in wastewater treatment  
 607 plants (WWTP) with a focus on integrating data assimilation techniques.

608 The objective was to identify the most effective approach for model-based  
609 control, taking into account real-world constraints such as limited data avail-  
610 ability and sensor noise. Data assimilation proved to be crucial for simulta-  
611 neously estimating model parameters, latent variables (such as unmeasured  
612 species concentrations), and quantifying prediction uncertainty which was  
613 rarely done before in the water research literature. Different data assimi-  
614 lation methods, including the Kalman filter, particle filter, and conditional  
615 particle filter, were combined with parameter estimation algorithms such as  
616 EM and MCEM to achieve robust model calibration.

617 The simulation results showed that the non-parametric black-box model  
618 outperformed all other models in terms of predictive accuracy and uncer-  
619 tainty estimation. This finding highlights the effectiveness of machine learn-  
620 ing when integrated with data assimilation techniques to capture informa-  
621 tion from the training set, even in the presence of limited data. The linear  
622 model, although quick to train, failed to capture the complex dynamics of  
623 the reactor. The white-box model, while potentially powerful, proved diffi-  
624 cult to train, time-consuming, and highly dependent on good initial model  
625 parameters, making it less generalizable. The grey-box model offered a good  
626 balance between accuracy and computational efficiency, though certain sim-  
627 plifications may have hindered its performance.

628 Interestingly, the addition of an extra sensor, such as an oxygen sensor,  
629 did not improve the model's performance when the target species (ammo-  
630 nium) was already being observed. This finding suggests that, at least for  
631 the non-parametric black-box model, investing in additional sensors may not  
632 be justified if the primary species of interest is already being monitored.

633 Experiments on the real system showed that the non-parametric black-  
634 box model was able to capture the general dynamics of ammonium concentra-  
635 tion in a real wastewater treatment plant. However, the model's performance  
636 was somewhat reduced compared to the simulation results, likely due to the  
637 variability in the input concentration, which was not accounted for in the  
638 model.

639 In conclusion, this study highlights the importance of data assimilation for  
640 robust WWTP modeling, especially when data are limited. It also demon-  
641 strates the potential of non-parametric "black-box" models for prediction  
642 and uncertainty quantification in this context. The study also highlighted  
643 the importance of accounting for input data variability, such as inlet concen-  
644 trations, to further improve model accuracy. Future research should focus on  
645 incorporating additional covariates, such as time of day, precipitation, and  
646 other environmental factors, to enhance model accuracy under real-world  
647 conditions. Investigating the impact of noise levels, noise distribution, and  
648 observation time steps could also provide valuable insights for model re-  
649 finement. Further improvements in confidence interval estimation could be  
650 achieved by incorporating heteroscedastic noise, allowing for more accurate  
651 representation of uncertainty in model predictions (Cheng et al., 2023). Ad-  
652 ditionally, future research could investigate the applicability of the proposed  
653 approach to WWTPs with more complex configurations, such as plants with  
654 parallel or serial tanks, or those receiving highly concentrated influents (e.g.,  
655 reject water from sludge treatment), to assess its robustness under varying  
656 operational conditions.

657 **Declaration of Competing Interest**

658 The authors declare that they have no known competing financial inter-  
659 ests or personal relationships that could have appeared to influence the work  
660 reported in this paper.

661 **Acknowledgments**

662 This material is based upon work supported by the ANRT (Association  
663 nationale de la recherche et de la technologie) with a CIFRE fellowship  
664 [grant numbers 2022/1582] in collaboration between the IRMAR (Institut  
665 de Recherche Mathématique de Rennes), the SATIE (Systèmes et applica-  
666 tions des technologies de l’information et de l’énergie) and the Purecontrol  
667 Company.

668 **References**

- 669 RTE, Futurs énergétiques 2050, 2021.
- 670 D. J. Reardon, Turning down the power, *Civil Engineering* 65 (1995) 54.
- 671 B. Chachuat, N. Roche, M. A. Latifi, Long-term optimal aeration strate-  
672 gies for small-size alternating activated sludge treatment plants, *Chemical*  
673 *Engineering and Processing: Process Intensification* 44 (2005) 591–604.
- 674 S. Balku, Comparison between alternating aerobic–anoxic and conventional  
675 activated sludge systems, *Water research* 41 (2007) 2220–2228.

- 676 L. J. S. Lukasse, K. J. Keesman, G. van Straten, A recursively identified  
677 model for short-term predictions of  $\text{nh}_4/\text{no}_3$  – concentrations in alternating  
678 activated sludge processes, *Journal of Process Control* 9 (1999) 87–100.
- 679 G. Fu, Y. Jin, S. Sun, Z. Yuan, D. Butler, The role of deep learning in urban  
680 water management: A critical review, *Water Research* 223 (2022) 118973.
- 681 R. S. Sutton, A. G. Barto, Reinforcement Learning, Adaptive Computation  
682 and Machine Learning series, 2 ed., Bradford Books, Cambridge, MA,  
683 2018.
- 684 G. Dulac-Arnold, D. J. Mankowitz, T. Hester, Challenges of real-world rein-  
685 forcement learning, *CoRR* abs/1904.12901 (2019).
- 686 M. Andreides, P. Dolejš, J. Bartáček, The prediction of wwtp influent char-  
687 acteristics: Good practices and challenges, *Journal of Water Process En-  
688 gineering* 49 (2022) 103009.
- 689 E. E. Holmes, Derivation of an em algorithm for constrained and uncon-  
690 strained multivariate autoregressive state-space (marss) models, 2013.
- 691 N. Kantas, A. Doucet, S. S. Singh, J. Maciejowski, N. Chopin, On parti-  
692 cle methods for parameter estimation in state-space models, *Statistical  
693 Science* 30 (2015).
- 694 M. Henze, W. Gujer, T. Mino, M. van Loosdrecht, Activated sludge models  
695 ASM1, ASM2, ASM2d and ASM3, IWA Publishing, 2006.
- 696 B. Chachuat, *Methodologie d’Optimisation Dynamique et de Commande Op-*

697 timale des Petites Stations d’Épuration a Boues Actives, Ph.D. thesis,  
698 Institut National Polytechnique de Lorraine (INPL), 2001.

699 B. Boulkroune, Estimation de l’état des systèmes non linéaires à temps dis-  
700 cret. Application à une station d’épuration, Ph.D. thesis, Université Henri  
701 Poincaré - Nancy I, 2008.

702 K. V. Germaey, M. C. van Loosdrecht, M. Henze, M. Lind, S. B. Jørgensen,  
703 Activated sludge wastewater treatment plant modelling and simulation:  
704 state of the art, *Environmental Modelling and Software* 19 (2004) 763–  
705 783.

706 A. Sánchez, M. Katebi, M. Johnson, Optimal control of an alternating  
707 aerobic-anoxic wastewater treatment plant, *IFAC Proceedings Volumes*  
708 35 (2002) 411–416.

709 M. Alvi, D. Batstone, C. K. Mbamba, P. Keymer, T. French, A. Ward,  
710 J. Dwyer, R. Cardell-Oliver, Deep learning in wastewater treatment: a  
711 critical review, *Water Research* 245 (2023) 120518.

712 M. Alvi, T. French, R. Cardell-Oliver, P. Keymer, A. Ward, Cost effective  
713 soft sensing for wastewater treatment facilities, *IEEE Access* 10 (2022)  
714 55694–55708.

715 M. Bocquet, Surrogate modeling for the climate sciences dynamics with  
716 machine learning and data assimilation, *Frontiers in Applied Mathematics*  
717 and Statistics 9 (2023).

718 M. Bonavita, On some limitations of current machine learning weather pre-  
719 diction models, *Geophysical Research Letters* 51 (2024) e2023GL107377.

- 720 R. Lam, A. Sanchez-Gonzalez, M. Willson, P. Wirnsberger, M. Fortunato,  
721 F. Alet, S. Ravuri, T. Ewalds, Z. Eaton-Rosen, W. Hu, A. Merose,  
722 S. Hoyer, G. Holland, O. Vinyals, J. Stott, A. Pritzel, S. Mohamed,  
723 P. Battaglia, Learning skillful medium-range global weather forecasting,  
724 Science 382 (2023) 1416–1421.
- 725 R. Lguensat, Les nouveaux modèles de prévision météorologique basés sur  
726 l’intelligence artificielle: opportunité ou menace ?, La Météorologie (2023)  
727 011.
- 728 L. Verhaeghe, J. Verwaeren, G. Kirim, S. Daneshgar, P. A. Vanrollegheem,  
729 E. Torfs, Towards good modelling practice for parallel hybrid models for  
730 wastewater treatment processes, Water Science and Technology 89 (2024)  
731 2971–2990.
- 732 T. T. T. Chau, P. Ailliot, V. Monbet, An algorithm for non-parametric esti-  
733 mation in state–space models, Computational Statistics & Data Analysis  
734 153 (2021) 107062.
- 735 C.-S. Gómez-Quintero, I. Queinnec, State and disturbance estimation for an  
736 alternating activated sludge process, IFAC Proceedings Volumes 35 (2002)  
737 447–452. 15th IFAC World Congress.
- 738 P. A. Stentoft, T. Munk-Nielsen, L. Vezzaro, H. Madsen, P. S. Mikkelsen,  
739 J. K. Møller, Towards model predictive control: online predictions of  
740 ammonium and nitrate removal by using a stochastic asm, Water Science  
741 and Technology 79 (2018) 51–62.

- 742 A. Stare, N. Hvala, D. Vrečko, Modeling, identification, and validation of  
743 models for predictive ammonia control in a wastewater treatment plant—a  
744 case study, *ISA Transactions* 45 (2006) 159–174.
- 745 C. Yang, P. Seiler, E. Belia, G. T. Daigger, An adaptive real-time grey-box  
746 model for advanced control and operations in wrrfs, *Water Science and*  
747 *Technology* 84 (2021) 2353–2365.
- 748 R. Fablet, B. Chapron, L. Drumetz, E. Mémin, O. Pannekoucke, F. Rousseau,  
749 Learning variational data assimilation models and solvers, *Journal of Ad-*  
750 *vances in Modeling Earth Systems* 13 (2021) e2021MS002572.
- 751 M. A. Sacco, M. Pulido, J. J. Ruiz, P. Tandeo, On-line machine-learning  
752 forecast uncertainty estimation for sequential data assimilation, *Quarterly*  
753 *Journal of the Royal Meteorological Society* n/a (2024).
- 754 J. Durbin, S. J. Koopman, Time series analysis by state space methods, Ox-  
755 ford statistical science series, 2nd ed ed., Oxford University Press, Oxford,  
756 2012.
- 757 S. Cheng, C. Quilodrán-Casas, S. Ouala, A. Farchi, C. Liu, P. Tandeo, R. Fa-  
758 blet, D. Lucor, B. Iooss, J. Brajard, et al., Machine learning with data  
759 assimilation and uncertainty quantification for dynamical systems: a re-  
760 view, *IEEE/CAA Journal of Automatica Sinica* 10 (2023) 1361–1387.
- 761 M. Henze, M. C. M. V. Loosdrecht, G. A. Ekama, D. Brdjanovic, *Biological*  
762 *Wastewater Treatment*, IWA Publishing, 2008.
- 763 Metcalf & Eddy Inc., G. Tchobanoglous, F. L. Burton, R. Tsuchihashi, H. D.

- 764 Stensel, Wastewater engineering: Treatment and resource recovery, 5 ed.,  
765 McGraw-Hill Professional, 2013.
- 766 V. Bertret, R. L. G. Latimier, V. Monbet, A stochastic expectation maxi-  
767 mization algorithm for the estimation of wastewater treatment plant am-  
768 monium concentration, in: 2024 European Control Conference (ECC),  
769 2024, pp. 3527–3532.
- 770 R. Lguensat, P. Tandeo, P. Ailliot, M. Pulido, R. Fablet, The analog data  
771 assimilation, Monthly Weather Review 145 (2017) 4093 – 4107.
- 772 F. Stulp, O. Sigaud, Many regression algorithms, one unified model: A  
773 review, Neural Networks 69 (2015) 60–79.
- 774 D. A. Otchere, T. O. Arbi Ganat, R. Gholami, S. Ridha, Application of su-  
775 pervised machine learning paradigms in the prediction of petroleum reser-  
776 voir properties: Comparative analysis of ann and svm models, Journal of  
777 Petroleum Science and Engineering 200 (2021) 108182.
- 778 A. Carrassi, M. Bocquet, L. Bertino, G. Evensen, Data assimilation in the  
779 geosciences: An overview of methods, issues, and perspectives, WIREs  
780 Climate Change 9 (2018) e535.
- 781 Y. Sasaki, Some basic formalisms in numerical variational analysis, Monthly  
782 Weather Review 98 (1970) 875 – 883.
- 783 B. D. Anderson, J. B. Moore, Optimal filtering, Courier Corporation, 2012.
- 784 G. Evensen, P. J. v. Leeuwen, An ensemble kalman smoother for nonlinear  
785 dynamics, Monthly Weather Review 128 (2000) 1852–1867.

- 786 A. Doucet, A. M. Johansen, A tutorial on particle filtering and smooth-  
787 ing : fifteen years later, in: D. Crisan, B. Rozovskii (Eds.), The Oxford  
788 handbook of nonlinear filtering, Oxford handbooks in mathematics, Ox-  
789 ford University Press, Oxford ; N.Y., 2011, pp. 656–705.
- 790 D. Dreano, P. Tandeo, M. Pulido, B. Ait-El-Fquih, T. Chonavel, I. Hoteit,  
791 Estimating model-error covariances in nonlinear state-space models using  
792 kalman smoothing and the expectation–maximization algorithm, Quar-  
793 terly Journal of the Royal Meteorological Society 143 (2017) 1877–1885.
- 794 A. Doucet, N. De Freitas, N. J. Gordon, et al., Sequential Monte Carlo  
795 methods in practice, volume 1, Springer, 2001.
- 796 F. Lindsten, Particle filters and Markov chains for learning of dynamical  
797 systems, Ph.D. thesis, Linköping University Electronic Press, 2013.
- 798 F. Lindsten, M. I. Jordan, T. B. Schon, Particle gibbs with ancestor sampling,  
799 Journal of Machine Learning Research 15 (2014) 2145–2184.
- 800 F. Lindsten, An efficient stochastic approximation em algorithm using condi-  
801 tional particle filters, in: 2013 IEEE International Conference on Acoustics,  
802 Speech and Signal Processing, 2013, pp. 6274–6278.
- 803 A. P. Dempster, N. M. Laird, D. B. Rubin, Maximum likelihood from incom-  
804 plete data via the em algorithm, Journal of the Royal Statistical Society.  
805 Series B (Methodological) 39 (1977) 1–38.
- 806 G. Celeux, D. Chauveau, J. Diebolt, On Stochastic Versions of the EM Al-  
807 gorithm, Research Report RR-2514, INRIA, 1995.

- 808 L. Ljung, System Identification : Theory for the User, 2nd ed., Pearson  
809 Education, 1998.
- 810 E. Paul, S. Plisson-Saune, M. Mauret, J. Cantet, Process state evaluation  
811 of alternating oxic-anoxic activated sludge using orp, ph and do, Water  
812 Science and Technology 38 (1998) 299–306.
- 813 K. V. Gernaey, U. Jeppsson, P. A. Vanrolleghem, J. B. Copp, Benchmarking  
814 of Control Strategies for Wastewater Treatment Plants, IWA Publishing,  
815 2014.