

Identification of Biokinetic Models using the Concept of Extents

Alma Mašić,[†] Sriniketh Srinivasan,[‡] Julien Billeter,[‡] Dominique Bonvin,[‡] and
Kris Villez^{*,†}

[†]*Eawag: Swiss Federal Institute of Aquatic Science and Technology, Überlandstrasse 133,
CH-8600 Dübendorf, Switzerland*

[‡]*Laboratoire d'Automatique, École Polytechnique Fédérale de Lausanne, CH-1015
Lausanne, Switzerland*

E-mail: Kris.Villez@eawag.ch

1

2

*Manuscript and Supporting Information accepted for publication in
Environmental Science & Technology*

3

4

5

Abstract

6

7

8

9

10

11

12

13

The development of a wide array of process technologies to enable the shift from conventional biological wastewater treatment processes to resource recovery systems is matched by an increasing demand for predictive capabilities. Mathematical models are excellent tools to meet this demand. However, obtaining reliable and fit-for-purpose models remains a cumbersome task due to the inherent complexity of biological wastewater treatment processes. In this work, we present a first study in the context of environmental biotechnology that adopts and explores the use of extents as a way to simplify and streamline the dynamic process modeling task. In addition, the

14 extent-based modeling strategy is enhanced by optimal accounting for nonlinear alge-
15 braic equilibria and nonlinear measurement equations. Finally, a thorough discussion
16 of our results explains the benefits of extent-based modeling and its potential to turn
17 environmental process modeling into a highly automated task.

18 **Introduction**

19 Dynamic models are increasingly used to better understand, design, and operate environ-
20 mental processes.^{1,2} For biological wastewater treatment processes, the available activated
21 sludge model family^{3,4} has been used widely despite reported challenges in model identifi-
22 cation. These challenges relate to *(i)* the information content and the quality of calibration
23 data that limit practical identifiability,⁵⁻¹⁰ *(ii)* the lack of mechanistic understanding,^{11,12}
24 and *(iii)* nonlinear and non-convex properties.¹³⁻¹⁵ These issues are even more severe in the
25 case of decentralized treatment processes that are proposed to address fast societal dynamics
26 by providing straightforward upscaling of wastewater treatment operations.¹⁶ In addition,
27 both economical and political motives are driving a paradigm shift in objectives from en-
28 vironmental protection to a need to generate added-value products from wastewaters. To
29 ensure both product quality and economically optimal operation, resource recovery from
30 wastewater requires tight management and control of the involved processes. The urine
31 nitrification process for fertilizer production developed at Eawag is an example of this.¹⁷
32 Advanced control of such high-rate processes is not possible without detailed process un-
33 derstanding and predictive power. In addition, the diversity of the available technologies
34 is rapidly increasing. For this reason, fast development of reliable models is paramount to
35 attain sustainable urban water cycles.

36 In the past, model complexity has been tackled by means of model identification
37 protocols. Examples include *(i)* protocols that split model identification into steps corre-
38 sponding to major fractions of the medium¹⁸ and *(ii)* protocols based on iterative model
39 building.¹⁹ Despite these efforts, the aforementioned model identification challenges have

40 only been partly addressed. In this work, we focus on the development of a method that
41 deals with the nonlinear and non-convex nature of kinetic identification in biological process
42 modeling. In previous work,²⁰ a deterministic optimization method was found well suited
43 to estimate parameters in a simple model for biological nitrite oxidation. This optimization
44 method led to globally optimal parameter estimates. The same study demonstrated that a
45 standard approach based on gradient-based optimization fails to find good parameter esti-
46 mates. Unfortunately, deterministic global optimization is cumbersome when the number of
47 parameters is large.

48 To deal with the model structure selection and parameter estimation challenge, we adopt
49 an extent-based framework^{21–24} to enable the application of deterministic optimization meth-
50 ods to biological process models involving multiple reactions. The concept of extents allows
51 the transformation multivariate time series into a set of individual time series, each one re-
52 flecting the progress of a single reaction. This, in turn, enables the individual identification
53 of the rate law and the corresponding parameters for each of the biological reactions. In ad-
54 dition, the use of extents facilitates model diagnosis. The proposed extent-based modeling
55 methodology is demonstrated and benchmarked against a conventional approach by means
56 of a simulated experiment with a urine nitrification process model.²⁵ All symbols used in
57 this text are listed in Table 1.

58 Other factors complicating model identification include *(i)* the stochastic nature of en-
59 vironmental processes and *(ii)* the significant lack of identifiability of model structures and
60 parameters, further leading to significant uncertainty and correlated parameter estimates.
61 These issues are certainly important but not studied in this work. Instead, we focus on
62 solving model identification problems to global optimality given experimental data. This
63 also means that we assume that a proper experimental design has been executed.

64 Methods

65 Definitions

66 Species and Components

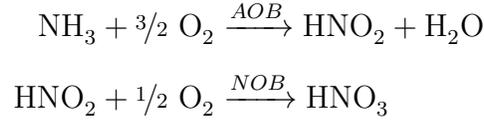
67 Consider a batch reactor with volume V containing S chemical species involved in R re-
68 actions. The numbers of moles are given as the S -dimensional vector \mathbf{n} . Among the R
69 reactions, R_k reactions are kinetically controlled, and R_e reactions are considered to be at
70 equilibrium, with $R = R_k + R_e$. The S species are split into S_k *kinetic species* that are only
71 involved in kinetically controlled reactions (i.e., *not* in equilibrium reactions) and S_e *equi-*
72 *librium species* that are involved in equilibrium (and possibly also in kinetically controlled)
73 reactions ($S = S_k + S_e$). The corresponding numbers of moles are \mathbf{n}_k and \mathbf{n}_e . *Equilibrium*
74 *components* are defined as the S_e molecular constituents that are involved in equilibrium
75 reactions and whose concentrations are conserved.²⁶ The $\bar{S} = S_k + S_e$ numbers of moles of
76 the kinetic species \mathbf{n}_k and the equilibrium components \mathbf{n}_c can be written as:

$$\bar{\mathbf{n}} = \begin{bmatrix} \mathbf{n}_k \\ \mathbf{n}_c \end{bmatrix} = \bar{\mathbf{E}} \mathbf{n} \quad (1)$$

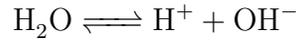
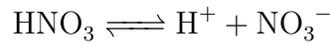
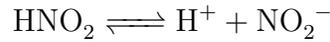
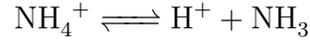
77 where $\bar{\mathbf{E}}$ of dimension $\bar{S} \times S$ relates the numbers of moles of all species \mathbf{n} to those of the
78 kinetic species and equilibrium components $\bar{\mathbf{n}}$.

Example. Let us illustrate the notations through a simplified urine nitrification process model²⁵ that is used in this work to simulate experimental data. This model is selected because it is an excellent example of a biological process model based on the activated sludge model family and involving rate-controlling acid-base equilibria. There are $R = 6$ reactions involving $S = 10$ species dissolved in water. The kinetically controlled reactions are the biological nitrification and nitrification by ammonia oxidizing bacteria (AOB) and

nitrite oxidizing bacteria (NOB), respectively, that is, $R_k = 2$:



The remaining reactions consist of $R_e = 4$ instantaneous acid-base equilibrium reactions:



79 The net growth of bacteria is assumed negligible. According to this reaction scheme, the
 80 $S = 10$ species are oxygen, ammonium, ammonia, nitrous acid, nitrite, nitric acid, nitrate,
 81 proton ions, hydroxyl ions, and water. Oxygen is *only* involved in the kinetically controlled
 82 reactions ($S_k = 1$). The remaining species are equilibrium species ($S_e = 9$). The numbers of
 83 moles are computed from the concentrations as follows:

$$\begin{aligned} \mathbf{n} &= V \left[[\text{O}_2] \quad [\text{NH}_4^+] \quad [\text{NH}_3] \quad [\text{HNO}_2] \quad [\text{NO}_2^-] \quad [\text{HNO}_3] \quad [\text{NO}_3^-] \quad [\text{H}^+] \quad [\text{OH}^-] \quad [\text{H}_2\text{O}] \right]^T \\ &= \begin{bmatrix} \mathbf{n}_k \\ \mathbf{n}_e \end{bmatrix} = \begin{bmatrix} V [\text{O}_2] \\ \mathbf{n}_e \end{bmatrix}. \end{aligned} \quad (2)$$

84 The $S_c = 5$ molecular constituents that are conserved in the equilibrium reactions are
 85 total ammonia (total ammonia nitrogen, TAN), total nitrite (TNO2), total nitrate (TNO3),
 86 total proton (TH), and total hydroxyl (TOH). With $S_k = 1$ (oxygen), the 6×10 matrix $\bar{\mathbf{E}}$
 87 reads:

$$\bar{\mathbf{E}} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 & 0 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 \end{bmatrix}. \quad (3)$$

88 **Balance Equations**

89 For a batch reactor, the differential mole balance equations are written as

$$\dot{\mathbf{n}}(t) = V \mathbf{N}^T \mathbf{r}(\mathbf{n}(t)/V), \quad \mathbf{n}(0) = \mathbf{n}_0 \quad (4)$$

90 with \mathbf{N} the $R \times S$ stoichiometric matrix, V the volume (assumed to be constant), \mathbf{r}
 91 the R -dimensional reaction rates, and \mathbf{n}_0 the S -dimensional initial numbers of moles. Upon
 92 pre-multiplying (4) by $\bar{\mathbf{E}}$, one obtains:

$$\dot{\bar{\mathbf{n}}}(t) = \bar{\mathbf{E}} \dot{\mathbf{n}}(t) = V \bar{\mathbf{E}} \mathbf{N}^T \mathbf{r}(\mathbf{n}(t)/V) = V \bar{\mathbf{N}}^T \mathbf{r}(\mathbf{n}(t)/V), \quad \bar{\mathbf{n}}(0) = \bar{\mathbf{n}}_0 \quad (5)$$

93 with $\bar{\mathbf{N}}$ the corresponding stoichiometric matrix of dimension $R \times \bar{S}$ and $\bar{\mathbf{n}}_0 = \bar{\mathbf{E}} \mathbf{n}_0$.
 94 Given $\bar{\mathbf{n}}$, the vector \mathbf{n} is obtained by solving the following system of $S = \bar{S} + R_e$ algebraic
 95 equations:²⁶

$$\bar{\mathbf{E}} \mathbf{n}(t) = \bar{\mathbf{n}}(t) \quad (6)$$

$$\mathbf{g}(\mathbf{n}(t)/V) = \mathbf{0}_{R_e} \quad (7)$$

96 where $\mathbf{g}(\cdot)$ expresses the R_e instantaneous equilibria. The dynamics of the component
 97 concentrations are functions of the kinetically controlled reactions only, that is, the rows
 98 of $\bar{\mathbf{N}}$ corresponding to the equilibrium reactions contain only zeros.²⁶ Hence, a reduced
 99 stoichiometric matrix $\bar{\mathbf{N}}_k$ can be defined as the matrix consisting of the rows of $\bar{\mathbf{N}}$ with at
 100 least one non-zero element. Following this, (5) reduces to:

$$\dot{\bar{\mathbf{n}}}(t) = V \bar{\mathbf{N}}_k^T \mathbf{r}_k(\mathbf{n}(t)/V), \quad \bar{\mathbf{n}}(0) = \bar{\mathbf{n}}_0 \quad (8)$$

101 with \mathbf{r}_k the kinetically controlled reaction rates.

102 **Example.** Following the aforementioned definitions, the stoichiometric matrix for all re-
 103 actions is

$$\mathbf{N} = \begin{bmatrix} -3/2 & 0 & -1 & 1 & 0 & 0 & 0 & 0 & 0 & 1 \\ -1/2 & 0 & 0 & -1 & 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & -1 & 1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & -1 & 1 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & -1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & -1 \end{bmatrix} \quad (9)$$

104 with the reduced stoichiometric matrix

$$\bar{\mathbf{N}}_k = \begin{bmatrix} -3/2 & -1 & 1 & 0 & 2 & 1 \\ -1/2 & 0 & -1 & 1 & 0 & 0 \end{bmatrix} \quad (10)$$

105 describing the $R_k = 2$ kinetically controlled reactions in terms of the $S_k = 1$ kinetic
 106 species and the $S_c = 5$ equilibrium components.

107 The rate laws for the biological oxidation reactions are:

$$\mathbf{r}_k = \begin{bmatrix} r_{\text{AOB}} \\ r_{\text{NOB}} \end{bmatrix} = \begin{bmatrix} [\text{NH}_3] / \left(\theta_{\text{AOB},1} + \theta_{\text{AOB},2} [\text{NH}_3] + \theta_{\text{AOB},3} [\text{NH}_3]^2 \right) \\ [\text{HNO}_2] / \left(\theta_{\text{NOB},1} + \theta_{\text{NOB},2} [\text{HNO}_2] \right) \end{bmatrix} \quad (11)$$

108 with the kinetic parameters $\theta_{\text{AOB},1}$, $\theta_{\text{AOB},2}$, $\theta_{\text{AOB},3}$, $\theta_{\text{NOB},1}$, and $\theta_{\text{NOB},2}$. The time depen-
 109 dence of rates and concentrations is omitted for the sake of conciseness. The two kinetic
 110 expressions correspond to Haldane and Monod kinetics, respectively. Since we assume that
 111 oxygen is sufficient for both oxidation processes, rate-limiting effects of oxygen can be safely
 112 ignored. The balance equations describing the equilibria cover four acid-base reactions so
 113 that (7) is

$$\mathbf{g}(\mathbf{n}/V) = \begin{bmatrix} \left([\text{H}^+] [\text{NH}_3] \right) / [\text{NH}_4^+] - 10^{-pK_{a,\text{NH}_4^+}} \\ \left([\text{H}^+] [\text{NO}_2^-] \right) / [\text{HNO}_2] - 10^{-pK_{a,\text{HNO}_2}} \\ \left([\text{H}^+] [\text{NO}_3^-] \right) / [\text{HNO}_3] - 10^{-pK_{a,\text{HNO}_3}} \\ [\text{H}^+] [\text{OH}^-] - 10^{-pK_w} \end{bmatrix} = \mathbf{0}_4. \quad (12)$$

114 The initial numbers of moles are $\mathbf{n}_{k,0} = V [c_{\text{O}_2,0} \ c_{\text{TAN},0} \ 0 \ 0 \ c_{\text{TH},0} \ c_{\text{TOH},0}]^T$. As
 115 there is no liquid entering or leaving the reactor during the reaction, the simulated data
 116 correspond to a typical batch test²⁷⁻³⁰ with a single pulse of ammonia dosed at the start

117 of the experiment. The initial concentrations of proton and hydroxyl component ($c_{\text{TH},0}$
 118 and $c_{\text{TOH},0}$) are set to values that satisfy the equilibrium equations and deliver a zero ion
 119 balance. The initial oxygen and water concentrations can be set arbitrarily and do not affect
 120 the reaction rates nor the equilibria.

121 Measurement Equations

122 During the simulated batch experiment, M measurements are obtained at H distinct time
 123 instants t_h , with $h = 1, \dots, H$ and $t_1 = 0$, as:

$$\tilde{\mathbf{y}}(t_h) = \mathbf{y}(t_h) + \boldsymbol{\epsilon}(t_h) = \mathbf{f}(\mathbf{n}(t_h)/V) + \boldsymbol{\epsilon}(t_h), \quad \boldsymbol{\epsilon}(t_h) \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_h) \quad (13)$$

124 with $\tilde{\mathbf{y}}(t_h)$ the M -dimensional vector of measurements, and $\mathbf{y}(t_h)$ the noise-free measured
 125 variables at time t_h . In words, the measurements are nonlinear functions of the species
 126 concentrations and are subject to additive Gaussian noise. The functions $\mathbf{f}(\cdot)$ are assumed
 127 continuous and differentiable.

128 **Example.** Measurements of the total ammonia, total nitrite and total nitrate concentra-
 129 tions and of pH are obtained. The noise-free measurements (13) are given as:

$$\mathbf{y} = \begin{bmatrix} y_{\text{TAN}} \\ y_{\text{TNO}_2} \\ y_{\text{TNO}_3} \\ y_{\text{pH}} \end{bmatrix} = \begin{bmatrix} [\text{NH}_4^+] + [\text{NH}_3] \\ [\text{HNO}_2] + [\text{NO}_2^-] \\ [\text{HNO}_3] + [\text{NO}_3^-] \\ -\log_{10}([\text{H}^+]) \end{bmatrix} = \begin{bmatrix} \mathbf{G} \mathbf{n}/V \\ -\log_{10}([\text{H}^+]) \end{bmatrix} \quad (14)$$

130 where \mathbf{G} is the measurement matrix

$$\mathbf{G} = \begin{bmatrix} 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 \end{bmatrix}. \quad (15)$$

131 Clearly, the first three measurements are linear in the species concentrations. In contrast,
 132 the pH measurement depends nonlinearly on the proton concentration, which in turn depends
 133 nonlinearly on the component concentrations via the algebraic equilibrium relationships in
 134 (7). In our simulated experiment, the measurement error variance-covariance matrix is
 135 considered time-invariant and defined as follows:

$$\Sigma_h = \Sigma = \text{diag} \left(\left[\sigma_{TAN}^2 \quad \sigma_{TNO2}^2 \quad \sigma_{TNO3}^2 \quad \sigma_{pH}^2 \right]^T \right), \quad h = 1, \dots, H \quad (16)$$

136 where the $\text{diag}(\cdot)$ operator creates a diagonal matrix from a column vector argument.

137 Data Generation and Problem Formulation

138 Data Generation

139 The main objective of this paper is to compare a new method for model identification to
 140 a more conventional approach. To properly compare the two methods, simulated measure-
 141 ments are used. These measurements are obtained by solving the DAE system ((6)-(8))
 142 with (3) and (10)-(12) from $t_1 = 0$ to $t_H = 10$ h. Measurements are obtained by means of
 143 (13)-(16) at regular intervals of 10 minutes so that $H = 61$. All parameter values used for
 144 simulation are given in Table 1.

145 **Problem Formulation**

146 The model identification problem consists in finding an appropriate model based on the
147 measurements from a pulse experiment. For each kinetically controlled reaction, a set of five
148 candidate rate laws are proposed. These are the zeroth-order, first-order, Monod, Tessier,
149 and Haldane rate laws given in Table 2. The initial conditions $\bar{\mathbf{n}}_0$, the stoichiometric matrix
150 \mathbf{N} , the equilibrium equations $\mathbf{g}(\cdot)$, the measurement equations $\mathbf{f}(\cdot)$, and the measurement
151 error variance-covariance matrices Σ_h are assumed to be known. Hence, the aim is therefore
152 to identify which of the candidate rate laws are appropriate for the two reactions, while
153 also estimating the corresponding kinetic parameters. In this work, feasible values for the
154 parameters are considered to be in the interval $[10^{-6}, 10^2]$.

155 **Notation.** The j^{th} candidate rate law for the i^{th} kinetically controlled reaction is referred
156 to as $r_{k,i}^{(j)}$. The corresponding parameter vectors are $\theta_i^{(j)}$. The number of candidate rate laws
157 for the i^{th} reaction is J_i , so that $j = 1, 2, \dots, J_i$. For a given choice of rate laws for the
158 kinetically controlled reactions, the parameter vector composed of the joint set of parameter
159 vectors for all reactions is denoted as Θ .

160 **Method 1: Simultaneous Model Identification**

161 The simultaneous model identification procedure is an exhaustive method that consists in
162 building a model for every possible combination of the candidate rate laws (\mathbf{r}_k) followed
163 by the estimation of all kinetic parameters (Θ) for each model. As indicated above, we
164 assume that the stoichiometry and equilibrium relations are known and the rate laws and
165 their parameters are to be identified. For a given selection of candidate rate laws, parame-
166 ter estimation is formulated mathematically as the following weighted least squares (WLS)
167 estimation problem:

$$\hat{\Theta} = \arg \min_{\Theta} \sum_{h=1}^H (\tilde{\mathbf{y}}(t_h) - \mathbf{y}(t_h))^T \Sigma_h^{-1} (\tilde{\mathbf{y}}(t_h) - \mathbf{y}(t_h)) \quad (17)$$

$$\text{s.t. } \mathbf{y}(t_h) = \mathbf{f}(\mathbf{n}(t_h)/V) \quad (18)$$

$$\mathbf{g}(\mathbf{n}(t)/V) = \mathbf{0}_{R_e} \quad (19)$$

$$\bar{\mathbf{E}} \mathbf{n}(t) = \bar{\mathbf{n}}(t) \quad (20)$$

$$\bar{\mathbf{n}}(t) = V \int_0^t \bar{\mathbf{N}}_k^T \mathbf{r}_k(\mathbf{n}(\tau)/V, \Theta) d\tau, \quad \bar{\mathbf{n}}(0) = \bar{\mathbf{n}}_0 \quad (21)$$

$$\Theta = \left[\boldsymbol{\theta}_1^{(j)T}, \dots, \boldsymbol{\theta}_i^{(j)T}, \dots, \boldsymbol{\theta}_{R_k}^{(j)T} \right]^T \quad (22)$$

168 During this estimation, the simulated system (18)-(22) is the same as the data-generating
 169 process, except for the rate laws in \mathbf{r}_k and the parameters therein. Because the measure-
 170 ment errors are assumed to be normally distributed according to (13), minimizing the WLS
 171 objective corresponds to a maximum-likelihood estimation (MLE).

172 The optimization problem (17)-(22) is solved by means of the Nelder-Mead simplex al-
 173 gorithm.³¹ This algorithm is initiated with parameter values at the center of the feasible
 174 intervals considered above. The total number of models whose parameters need to be es-
 175 timated equals the product of the numbers of candidate rate laws, $\prod_i J_i$. Following the
 176 parameter estimation for each of these models, a well-fitting model is selected from the com-
 177 plete set of models by trading off the WLS objective (17) against parsimony. To this end,
 178 the WLS objective is equivalently expressed as the weighted root mean squared residual
 179 (WRMSR):

$$WRMSR = \sqrt{\frac{1}{H \cdot M} \sum_{h=1}^H (\tilde{\mathbf{y}}(t_h) - \mathbf{y}(t_h))^T \Sigma_h^{-1} (\tilde{\mathbf{y}}(t_h) - \mathbf{y}(t_h))} \quad (23)$$

180 **Example.** To model the simulated process, five different candidate rate laws are consid-

181 ered for each of the two biological reactions ($J_1 = J_2 = 5$). The number of distinct models
182 whose parameters are estimated is therefore $\prod_i J_i = J_1 \cdot J_2 = 5 \cdot 5 = 25$.

183 **Method 2: Incremental Model Identification via Extents**

184 This subsection introduces the concept of extents of reaction and shows how to compute
185 them from the measured numbers of moles. The computed extents, named experimental
186 extents, are then used to identify the kinetics of each reaction individually, thereby making
187 the procedure incremental. Finally, the same measurements are used to fine-tune the kinetic
188 parameters for the global model.

189 **Definition of Extents**

190 In batch reactors, the extents of reaction $\mathbf{x}(t)$ can be defined by means of the following
191 integral:

$$\mathbf{n}(t) = \mathbf{n}_0 + V \int_0^t \mathbf{N}^T \mathbf{r}(\mathbf{n}(\tau)/V) d\tau = \mathbf{n}_0 + \mathbf{N}^T \mathbf{x}(t). \quad (24)$$

192 In words, an extent of reaction expresses the progress of the corresponding reaction in
193 terms of the numbers of moles of the product it has produced since $t = 0$. This definition
194 can be applied to multiphase systems as well.²¹ In what follows, unless mentioned otherwise,
195 the term extent refers specifically to the extent of a kinetically controlled reaction. Equation
196 (8) can be integrated to give:

$$\bar{\mathbf{n}}(t) = \bar{\mathbf{n}}_0 + V \int_0^t \bar{\mathbf{N}}_k^T \mathbf{r}_k(\mathbf{n}(\tau)/V) d\tau = \bar{\mathbf{n}}_0 + \bar{\mathbf{N}}_k^T \mathbf{x}_k(t). \quad (25)$$

197 Reformulating the balance equations (21)-(20) in terms of extents allows the selection of
198 rate laws and estimating parameters for each reaction individually. To do so, the available

199 measurements are first transformed into *experimental extents*. After this transformation, and
 200 for each reaction individually, selected rate laws can be fitted to the experimental extents.
 201 These steps are explained next.

202 Step 1 – Computation of Experimental Extents

203 The extents of reaction for the kinetically controlled reactions can be computed by solving
 204 the following WLS problem for each sampling instant t_h :

$$\tilde{\mathbf{x}}_k(t_h) = \arg \min_{\mathbf{x}_k(t_h)} (\tilde{\mathbf{y}}(t_h) - \mathbf{y}(t_h))^T \boldsymbol{\Sigma}_h^{-1} (\tilde{\mathbf{y}}(t_h) - \mathbf{y}(t_h)) \quad (26)$$

$$\text{s.t. } \mathbf{y}(t_h) = \mathbf{f}(\mathbf{n}(t_h)/V) \quad (27)$$

$$\mathbf{g}(\mathbf{n}(t_h)/V) = \mathbf{0}_{R_e} \quad (28)$$

$$\bar{\mathbf{E}} \mathbf{n}(t_h) = \bar{\mathbf{n}}_0 + \bar{\mathbf{N}}_k^T \mathbf{x}_k(t_h) \quad (29)$$

205 where (26) is the objective function, (27) expresses the expected measurements as func-
 206 tions of the numbers of moles of the species, (28) expresses the algebraic equilibria, and
 207 (29) relates the extents of the kinetically controlled reactions to the number of moles of the
 208 species. As above, minimizing the WLS objective to its global minimum corresponds to MLE.
 209 In general, the above problem is nonlinear, non-convex, and therefore solved numerically. In
 210 special cases, an analytic solution can be provided.²⁶

211 The initial numbers of moles $\bar{\mathbf{n}}_0$, and volume V are assumed to be known. Hence,
 212 one can compute the point-wise approximation $\boldsymbol{\Lambda}_h$ to the variance-covariance matrix of the
 213 experimental extents as the inverse of the Fisher information matrix $\mathbf{J}_h^T \boldsymbol{\Sigma}_h^{-1} \mathbf{J}_h$, where \mathbf{J}_h is
 214 the Jacobian matrix, with $\mathbf{J}_h(m, i) = \partial y_m / \partial x_{k,i} |_{\tilde{\mathbf{x}}_k(t_h)}$. The elements of \mathbf{J}_h are computed by
 215 numerical differentiation unless analytical derivatives are available. This procedure allows
 216 writing the following approximate distribution for the extent estimation errors (i.e., the
 217 difference between the experimental extents $\tilde{\mathbf{x}}_k$ and the true extents \mathbf{x}_k):

$$\tilde{\mathbf{x}}_k(t_h) - \mathbf{x}_k(t_h) \sim \mathcal{N}(\mathbf{0}_{R_k}, \mathbf{\Lambda}_h). \quad (30)$$

218 **Example.** In the simulated experiment, it follows from (1), (14), and (25) that the Jaco-
 219 bian consists of three rows that are computed analytically and a fourth row that is evaluated
 220 numerically:

$$\mathbf{J}_h = \begin{bmatrix} \frac{1}{V} \mathbf{G} \bar{\mathbf{E}}^+ \bar{\mathbf{N}}_k^T \\ \left. \frac{\partial pH}{\partial \tilde{\mathbf{x}}_k} \right|_{\tilde{\mathbf{x}}_k(t_h)} \end{bmatrix} \quad (31)$$

221 with the superscript $(\cdot)^+$ indicating the Moore-Penrose pseudo-inverse. The i^{th} element
 222 of the last row is computed as $\frac{pH_1 - pH_0}{\delta}$, with pH_0 and pH_1 the pH values obtained by solving
 223 (27)-(29) at $\tilde{\mathbf{x}}_k(t_h)$ and $\tilde{\mathbf{x}}_k(t_h) + \Delta_i$, with Δ_i a vector with the small number δ in its i^{th}
 224 position and zeros elsewhere.

225 Step 2 – Extent Modeling

226 The original identification problem (17)-(22) is now simplified by fitting the rate laws to
 227 the experimental extents instead of to the original measurements and by estimating the
 228 parameters of a single reaction at the time. The idea is to model each reaction by optimizing
 229 the fit to the corresponding experimental extent, $\tilde{x}_{k,i}$, with $i = 1, \dots, R_k$. However, since
 230 the reaction rate $r_{k,i}$ is a function of concentrations that might depend on the progress of
 231 several reactions, one estimates the contribution of the *other* reactions from measurements.²¹
 232 This results in the following optimization problem for the j th candidate rate law for the i th
 233 kinetically controlled reaction:

$$\hat{\boldsymbol{\theta}}_i^{(j)} = \arg \min_{\boldsymbol{\theta}_i^{(j)}} \text{ssq}_i := \sum_{h=1}^H \frac{(\tilde{x}_{k,i}(t_h) - x_{k,i}(t_h))^2}{\lambda_{i,h}} \quad (32)$$

$$\text{s.t. } \mathbf{g}(\mathbf{n}(t)/V) = \mathbf{0}_{R_e} \quad (33)$$

$$\bar{\mathbf{E}} \mathbf{n}(t) = \bar{\mathbf{n}}_0 + \bar{\mathbf{N}}_k^T \mathbf{x}_k(t) \quad (34)$$

$$\forall r = 1, \dots, R_k : \quad (35)$$

$$x_{k,r}(t) = \begin{cases} V \int_0^t r_{k,i}^{(j)} \left(\frac{\mathbf{n}(\tau)}{V}, \boldsymbol{\theta}_i^{(j)} \right) d\tau, & x_{k,i}(0) = 0 \quad \text{if } r = i \\ \mathcal{I}(\mathbf{t}, \tilde{\mathbf{x}}_{k,r}, t), & \text{if } r \neq i \end{cases}$$

234 where $\lambda_{i,h} := \boldsymbol{\Lambda}_h(i, i)$, $\mathbf{t} = [t_1, t_2, \dots, t_h, \dots, t_H]$ and with the operator $\mathcal{I}(\cdot)$ defined as

$$\forall t \in \{t : t_l \leq t \leq t_{l+1}\} : \quad \mathcal{I}(\mathbf{t}, \tilde{\mathbf{x}}_{k,r}, t) := \tilde{x}_{k,r}(t_l) + \left(\tilde{x}_{k,r}(t_{l+1}) - \tilde{x}_{k,r}(t_l) \right) \frac{t - t_l}{t_{l+1} - t_l}. \quad (36)$$

235 In the above problem, (32) is the objective function expressing that the i th predicted
 236 extent should be as close as possible to the corresponding experimental extents in the WLS
 237 sense. As before, (33) and (34) express the algebraic equilibria and the relationships be-
 238 tween the extents of the kinetically controlled reactions and the number of moles of all
 239 species. Equation (35) indicates that the predicted extents stem from (i) the simulated i th
 240 reaction, and (ii) piecewise linear interpolation of the experimental extents for the other ki-
 241 netically controlled reactions. The most important consequence of this method is that only
 242 the kinetic parameters of the i th candidate rate law appear in the optimization problem.
 243 Indeed, the interpolation of the experimental extents (36) implies that the kinetic param-
 244 eters of the corresponding reactions are not needed. The original optimization problem is
 245 thereby replaced by multiple optimization problems involving a univariate system including
 246 only one reaction. Furthermore, the modification also means that one does not need to know
 247 the structure of the rate laws corresponding to the interpolated experimental extents, that

248 is, the best candidate rate law for each reaction can be found independently of the rate laws
 249 for the other reactions.

250 The second method allows solving each individual parameter estimation problem to global
 251 optimality by means of the branch-and-bound algorithm proposed earlier.²⁰ This way, the
 252 best parameter values are guaranteed to be found within the considered feasible intervals.
 253 The bounding procedures required for this algorithm are given in the Supporting Information.
 254 With each candidate rate law and the associated optimal parameters $\hat{\theta}_i^{(j)}$, one obtains the
 255 modeled extent $\hat{x}_{k,i}^{(j)}$ and the following extent-specific WRMSR:

$$WRMSR_i^{(j)} = \sqrt{\frac{1}{H} \sum_{h=1}^H \frac{\left(\tilde{x}_{k,i}(t_h) - \hat{x}_{k,i}^{(j)}(t_h)\right)^2}{\lambda_{i,h}}}. \quad (37)$$

256 The rate law $\hat{r}_{k,i}$ is selected by trading off the WRMSR against parsimony. This is
 257 repeated for every reaction, which means that the number of parameter estimation problems
 258 to be solved now equals the sum of the numbers of candidate rate laws, $\sum_i J_i$. In addition,
 259 the number of parameters that are estimated in each problem is generally lower than the
 260 number of parameters estimated with the first method (17)-(22).

261 **Example.** With 5 candidate rate laws considered for each reaction, $\sum_i J_i = J_1 + J_2 = 10$
 262 instances of the parameter estimation problem need to be solved. The number of parameters
 263 that are estimated in each problem ranges from 1 (e.g. zeroth-order rate law) to 3 (Haldane).
 264 In comparison, the first method requires the estimation of 2 (zeroth-order rate law for both
 265 reactions) up to 6 (Haldane rate law for both reactions) parameters at once.

266 Step 3 – Model Fine-Tuning

267 Following the rate-law selection, the model parameters are fine-tuned by simultaneously
 268 estimating all kinetic parameters via (17)-(22). As in the first method, this is done using

269 the Nelder-Mead simplex algorithm. In contrast to the first method, this algorithm is now
270 executed for only one model containing the rate laws selected in Step 2 and is initiated with
271 the corresponding parameter estimates obtained in Step 2.

272 **Results**

273 **Process Simulation**

274 The nitrification model (6)-(15) is used to generate concentration and pH measurement
275 series. The results are shown in Fig. 1. One can see a fairly distinct separation in time of the
276 two reactions with the TNO2 concentration rising to 50% of the original TAN concentration
277 at about 4 h. Before (after) this time, a net production (consumption) of TNO2 is observed.
278 The figure also shows the free ammonia concentration $[\text{NH}_3]$. The ammonia oxidation stops
279 when this concentration reaches zero. The nitrite and nitrate ion concentrations are nearly
280 indistinguishable from the total nitrite and total nitrate concentrations (not shown). At the
281 end of the experiment, about half of the available TAN is converted via nitrite to nitrate.
282 The limited buffering capacity in the simulated system causes fairly large changes in pH.
283 Additive Gaussian noise is simulated added to generate realistic measurements.

284 **Method 1: Simultaneous Model Identification**

285 The kinetic parameters of 25 different models, each with a unique pair of rate laws for the
286 first and second reactions, are estimated by solving (17)-(22). The resulting WRMSR values
287 shown in Fig. 2 range from 6.57 to 37.86. These values indicate the model prediction error
288 standard deviation relative to the measurement error standard deviation. Assuming the
289 correct model, it exhibits a χ^2 -distribution with a mean value of 1 and a right-sided 99%
290 confidence limit of 1.11. The graph also shows the WRMSR value of 1.01 obtained with
291 the noisy measurements and the true model including its parameters. This WRMSR is very
292 close to the expected value of 1. Note that the best model gives a WRMSR value that is

293 6.57 times larger than the WRMSR value obtained with method 2 (see below). Clearly,
294 this method is unfit to find a good model. In all cases, including the case involving the
295 true rate laws used for simulation, only a locally optimal parameter set could be found. In
296 addition, the best model (Model 16) includes the Tessier rate law for the first reaction and
297 the zeroth-order rate law for the second reaction, which does not correspond to the true rate
298 laws. Additional results, including simulations using each of the 25 models after parameter
299 estimation, are included in the Supporting Information.

300 **Method 2: Incremental Model Identification via Extents**

301 **Step 1 – Extent Computation**

302 The extents computed by solving (26)-(29) using the TAN, TNO₂, TNO₃, and pH mea-
303 surements are shown in Fig. 3(a). The confidence bands for the experimental extents vary
304 with time, in particular for the first extent. High precision is obtained at the beginning
305 and during most of the second half of the experiment. However, during the first half, the
306 uncertainty first increases and then decreases. At the end of the experiment, the uncertainty
307 increases again. These effects are due to the nonlinear propagation of the pH measurement
308 error through the measurement and algebraic equations. The ellipsoidal confidence regions
309 at 0.5, 1.5, 2.5, 3.5, and 4.5 h are shown in Fig. 3(b). The orientation of the confidence
310 region becomes more oblique with increased uncertainty in the first extent.

311 **Step 2 – Extent Modeling**

312 The global solutions to (32)-(35), obtained for every reaction and every candidate rate law,
313 are discussed next.

314 **Modeling the First Extent.** The best fits of the first extent obtained with the various
315 candidate rate laws are shown in Fig. 4(a). It is clear that the zeroth- and first-order models
316 do not fit the experimental extents well. The Monod and Tessier models fit better, yet they

317 over-estimate the experimental extent. This is clearly visible in Fig. 4(b-c), where the model
318 errors are shown. In contrast, the Haldane rate law fits the extent profile well. As such,
319 the Haldane model is easily selected as the best among the model candidates. In Fig. 5, the
320 WRMSR values (37) are given with a 95% upper control limit based on the corresponding
321 χ^2 -statistics. Based on this statistic, all models except for the Haldane model are rejected
322 for the first extent.

323 **Modeling the Second Extent.** The best fits of the second extent are visualized in
324 Fig. 6(a). Here, all rate models fit the experimental extents reasonably well, except for
325 the zeroth-order model. The first-order model leads to visibly auto-correlated residuals
326 Fig. 6(b-c). This is also evident from the WRMSR values (Fig. 5), on the basis of which the
327 zeroth-order and first-order models are rejected. In this case, a parsimonious model is chosen
328 among the three remaining candidates. The Monod model delivers the best fit among the
329 simplest candidates (Monod and Tessier). An alternative approach may consist in designing
330 an experiment that enables better discrimination of the remaining rate laws. This is not
331 explored in this work.

332 **Step 3 – Model Fine-Tuning**

333 The model structure consisting of the two selected rate laws, namely, Haldane and Monod, is
334 used next to fine-tune the model parameter via the simultaneous approach (17)-(22). Fig. 7
335 compares the simulated concentration and pH values with the predictions of the identified
336 models prior and after fine-tuning. These three simulations are hard to distinguish from
337 each other. The resulting overall WRMSR (23) equals 1.0013 and is shown in Fig. 2. Most
338 importantly, the extent-based model identification procedure has delivered a well-fitting set
339 of rate laws and kinetic parameter estimates. Furthermore, the selected rate laws are exactly
340 those used to generate the simulated experimental measurements. The parameter estimates
341 deviate at most 10% from their true values, except $\theta_{AOB,2}$ which deviates by about 30%. Such

342 deviations are typical for biokinetic wastewater treatment models and are in part explained
343 by correlation between parameter estimates.

344 Discussion

345 The results presented above are now interpreted in a broader biokinetic modeling context.

346 **Interpretation of the Results.** In this study, the concept of extents is introduced for the
347 first time for the purpose of dynamic modeling of an environmental biochemical process. By
348 means of a simplified biokinetic model of the urine nitrification process and simulated batch
349 experiments, several benefits of the extent-based modeling approach have been demonstrated.
350 Concretely, the identification of biokinetic models via extents:

- 351 • allows using deterministic optimization methods to obtain excellent parameter esti-
352 mates. Despite the fact that the individual extent modeling steps only approximate
353 the original model identification problem, one can obtain a well-fitting model. Most im-
354 portantly, the convergence to local optima as observed with a conventional parameter
355 estimation method can be avoided.
- 356 • provides an intuitive diagnostic tool for modeling. Indeed, extent-modeling indicates
357 whether a reaction can be modeled appropriately with a given candidate rate law, thus
358 allowing modelers to pay more attention to reactions that are more difficult to model.
359 Similarly, this approach indicates whether sufficient information is available within a
360 given experimental data set to discriminate between candidate rate laws.
- 361 • reduces a model selection problem that is polynomial in the number of candidate
362 rate laws to a model selection problem that is linear in this number. In this study,
363 the extent-based modeling method required solving 10 parameter estimation problems
364 involving 1 to 3 parameters, whereas the conventional simultaneous approach required
365 solving 25 parameter estimation problems involving 2 to 6 parameters.

366 It is of special importance that the extent-based model identification method is the
367 only method delivering an acceptable model. Indeed, the conventional model identification
368 method did not result in an acceptable model, despite the apparent simplicity of the studied
369 process and simulated experiment.

370 **Links in Prior Work.** While the concept of extents is new in the context of dynamic
371 modeling of environmental processes, it is important to note that a number of important
372 concepts in use today are somewhat similar. For instance, the integral defined by the area
373 under the oxygen uptake rate curve, a.k.a. *respirogram*, is matched to the total accumulated
374 oxygen uptake in typical respirometric experiments.³² Similar concepts include *accumulated*
375 *methane production*³³ and *number of base pulses*.³⁴ It is also interesting to note that the
376 accumulated cellulose solubilisation has been described as the *extent of solubilisation*.³⁵
377 However, this is without links to the general concept of extents. The most important dif-
378 ference between extents and the concepts already in use is that extents reflect individual
379 processes rather than several simultaneous processes. So far, model reduction on the basis
380 of the concept of reaction invariants³⁶ is the only related application known in the environ-
381 mental engineering sciences. We expect tangible benefits from a broader and systematic use
382 of extents, including those mentioned above.

383 **Analysis of the Extent-based Modeling Method.** In the general case, the extent-
384 based modeling method does not solve the exact same problem as the conventional simul-
385 taneous modeling method. Extent-based modeling solves the same problem (17)-(22) only
386 if (i) the measurement equations are linear and there are no nonlinear algebraic equations
387 involved in the extent computations (26)-(29), (ii) the off-diagonal elements of the matrices
388 $\mathbf{\Lambda}_h$ are equal to zero, that is, in absence of correlation between experimental extents, and
389 (iii) the reaction rates can be expressed as functions of the modeled extents. These require-
390 ments are rarely satisfied so that the resulting parameter estimates likely deviate from those
391 obtained by solving (17)-(22). However, the extent-based modeling framework is particularly

392 useful when solving (17)-(22) to global optimality is difficult or computationally prohibitive.

393 In computing the solution to (26)-(29), one can encounter different situations:

- 394 1. The first situation occurs when the available measurements are linear in the extents
395 of the kinetically controlled reactions and do not depend on the equilibrium species
396 concentrations. In this case, one can discard all nonlinear (equilibrium) equations and
397 an analytic solution for the extents can be found.

- 398 2. The second situation occurs when the number of measured variables matches the num-
399 ber of computed extents exactly, thereby resulting in a fully determined system (hence
400 no need for optimization). In this situation, one can find extents that make the objec-
401 tive function (26) equal to zero, while satisfying (27)-(29). The solution can therefore
402 be obtained by solving the equation system (27)-(29) numerically. In the process con-
403 sidered in this work, this situation would occur if the pH and one of the remaining
404 variables (TAN, TNO2, TNO3) were measured (not demonstrated).

- 405 3. The third situation occurs when the number of measured variables exceeds the number
406 of computed extents (overdetermined system). This corresponds to the case studied
407 in this work (TAN, TNO2, TNO3, and pH measured). One approach consists of
408 discarding (26) and solving (27)-(29) in a least-squares sense.²⁶ When doing so, the
409 experimental extents are not the solution to (26)-(29). We recommend solving (26)-(29)
410 exactly, as in this work, to obtain experimental extents that are WLS-optimal.

411 In its current form, the proposed extent-based modeling method assumes a closed batch
412 process whose stoichiometric matrix and the algebraic equilibrium equations are known or
413 estimated precisely. However, this is not true in general. The method presented here can
414 easily be expanded to account for mass transfer as well as gas-liquid transfer as demonstrated
415 already.^{26,37} The main reason this has not been included here is to maintain a clear presen-
416 tation of the developed method. Not knowing the stoichiometric matrix or the algebraic
417 equilibrium equations means that, prior to modeling via extents, one may use target factor

418 analysis³⁸ to identify the stoichiometric matrix or detailed physico-chemical analysis to ob-
419 tain a model for acid-base and salt speciation. However, certain situations allow using the
420 extent-based modeling framework to estimate equilibrium parameters²⁶ as well as stoichio-
421 metric parameters.³⁹ Even more critical is the fact that extent-based model identification
422 requires at least as many measured variables as there are kinetically controlled reactions.
423 When this requirement is not met, one can opt to partition the model identification problem
424 into smaller problems which include more than one reaction.³⁹

425 **Methodological Improvements.** Methodologically speaking, this work adds four ele-
426 ments to the extent-based modeling framework, namely:

- 427 • Extent computation with measurements that are nonlinear in the species concentra-
428 tions.
- 429 • Optimal estimation of the experimental extents when (27)-(29) involves more measured
430 variables than extents, that is, in the overdetermined case.
- 431 • Accounting for nonlinear effects during experimental extent computation by means of
432 a Laplacian approximation of their distribution.
- 433 • Extent-based modeling and deterministic global optimization are combined for the first
434 time into a single model identification framework.

435 **Future Work.** The developments in this study are considered critical steps towards a first
436 real-world application of the extent-based modeling of environmental processes. However, the
437 following aspects call for further development and testing of the method prior to experimental
438 validation in full-scale wastewater treatment systems:

- 439 • Realistic sensor data. So far, measurement devices are considered to exhibit an in-
440 stantaneous response within the extent-based modeling framework. However, typical
441 devices respond dynamically to the measured variable.^{40,41} Explicit accounting of sen-
442 sor dynamics is not feasible yet in the extent-based modeling framework.

- 443 • Laboratory validation. Several aspects of real biological processes have been ignored to
444 facilitate the introduction of extent-based modeling. The ignored elements include *(i)*
445 bacterial growth and decay processes, *(ii)* complex composition of actual wastewater,
446 and *(iii)* complex physico-chemical reaction systems in high-strength wastewater such
447 as source-separated urine. The first element only affects the extent-based methodology
448 due to a lack of extent observability. This can be accounted for in special cases³⁹ but
449 may prove difficult to address in general.⁴² The second and third element affect both
450 modeling methods used in this study and are being addressed currently by adopting
451 a more realistic physico-chemical urine composition and associated reaction system in
452 view of a lab-scale validation.
- 453 • Prior knowledge. In this work, the reactor volume V , the initial conditions $\bar{\mathbf{n}}_0$, and the
454 stoichiometric matrix \mathbf{N} are considered known. Methods permitting the estimation of
455 these variables and parameters remain to be investigated.
- 456 • Completeness of the candidate rate laws. In this work, we have assumed that the set
457 of candidate rate laws includes the true rate laws in the data-generating process. This
458 is not true in general. An alternative model structure based on shape constrained
459 splines can address this problem.⁴³ So far, this type of models has only been applied to
460 monoculture processes. Its use in connection with the extent-based model identification
461 remains to be evaluated.

462 **Acknowledgement**

463 This study was made possible by Eawag Discretionary Funds (grant no.: 5221.00492.009.03,
464 project: DF2015/EMISSUN). All results were obtained by use of Matlab⁴⁴ and the Spike_O
465 toolbox for optimization.⁴⁵

Table 1: List of symbols and parameter values used for simulation. Values in parentheses refer to the best-available estimates.

Symbol	Description	Value	Unit
c_S	Substrate concentration	–	mol·L ⁻¹
$c_{\text{TAN},0}$	Initial TAN concentration	0.35	mol·L ⁻¹
$c_{\text{TH},0}, c_{\text{TOH},0}, c_{\text{O}_2,0}$	Initial concentrations	–	mol·L ⁻¹
$\bar{\mathbf{E}}$	Matrix defining the set of kinetic species and conserved molecular constituents	–	–
$\mathbf{f}(\cdot)$	Measurement expressions	–	–
\mathbf{G}	Measurement gain matrix	–	–
$\mathbf{g}(\cdot)$	Algebraic equilibrium expressions	–	–
H	Number of samples	61	–
h	Measurement sample index	–	–
i	Reaction index	–	–
\mathbf{J}_h	Jacobian matrix at the h th sample time	–	–
J_i	Number of rate law candidates for the i th reaction	–	–
j	Rate law candidate index	–	–
M	Number of measurements	–	–
m	Measured variable index	–	–
$\mathbf{N}, \bar{\mathbf{N}}, \bar{\mathbf{N}}_k$	Stoichiometric matrices	–	–
\mathbf{n}	Numbers of moles of all species	–	mol
$\bar{\mathbf{n}}$	Numbers of moles of kinetic species and conserved molecular constituents	–	mol
$\mathbf{n}_0, \mathbf{n}_{k,0}, \bar{\mathbf{n}}_0$	Initial numbers of moles	–	mol
\mathbf{n}_c	Numbers of moles of conserved molecular constituents	–	mol
\mathbf{n}_e	Numbers of moles of equilibrium species	–	mol
\mathbf{n}_k	Numbers of moles of kinetic species	–	mol
pK_{a,NH_4^+}	Logarithmic acid dissociation constant of NH ₄ ⁺	+9.24	–
pK_{a,HNO_2}	Logarithmic acid dissociation constant of HNO ₂	+3.25	–
pK_{a,HNO_3}	Logarithmic acid dissociation constant of HNO ₃	-1.40	–

Table 1: List of symbols and parameter values used for simulation. Values in parentheses refer to the best-available estimates.

Symbol	Description	Value	Unit
pK_w	Logarithmic water dissociation constant	14	—
R	Number of reactions	6	—
R_e	Number of equilibrium reactions	4	—
R_k	Number of kinetic reactions	2	—
r	Reaction index	—	—
\mathbf{r}	Reaction rates	—	$\text{mol}\cdot\text{L}^{-1}\cdot\text{h}^{-1}$
$\mathbf{r}_k (r_{k,i})$	Reaction rates of the kinetically controlled reactions (of the i th kinetically controlled reaction)	—	$\text{mol}\cdot\text{L}^{-1}\cdot\text{h}^{-1}$
$\hat{r}_{k,i}$	Selected rate law for the i th kinetically controlled reaction	—	—
S	Number of chemical species	10	—
\bar{S}	Number of kinetic species and conserved molecular constituents	6	—
S_c	Number of conserved molecular constituents	5	—
S_e	Number of equilibrium species	9	—
S_k	Number of kinetic species	1	—
t, t_h	Time (of measurement)	—	h
V	Volume	1	L
$WRMSR$	Weighted root mean squared residual	—	—
$WRMSR_i^{(j)}$	WRMSR of the i th reaction with the j th candidate rate law	—	—
$\mathbf{x}_k (x_{k,i}, x_{k,r})$	Extents of the kinetically controlled reactions (of the i/r th kinetically controlled reaction)	—	mol
$\tilde{\mathbf{x}}_k (\tilde{x}_{k,i}, \tilde{\mathbf{x}}_{k,r}, \tilde{x}_{k,r})$	Experimental extents of kinetically controlled reactions (of the i/r th kinetically controlled reaction)	—	mol
$\hat{x}_{k,i}^{(j)}$	Extent estimate for the i th kinetically controlled reaction with the j th rate law candidate	—	mol
\mathbf{y}	Measured variables	—	—
$\tilde{\mathbf{y}}$	Measurements	—	—

Table 1: List of symbols and parameter values used for simulation. Values in parentheses refer to the best-available estimates.

Symbol	Description	Value	Unit
$y_{TAN}, y_{TNO_2}, y_{TNO_3}, y_{pH}$	Measured variables	–	
$\Delta (\Delta_i)$	Perturbation vector (for the i th reaction)	–	mol
δ	Perturbation parameter	$1 \cdot 10^{-12}$	mol
Θ	Kinetic parameters for all rate laws	–	
$\hat{\Theta}$	Kinetic parameter estimates for all rate laws	–	
$\theta (\theta_i^{(j)})$	Kinetic parameters (for the j th rate law candidate of i th reaction)	–	
$\hat{\theta}_i^{(j)}$	Parameter estimates for j th kinetic rate law candidate for the i th kinetically controlled reaction	–	
$\theta_{AOB,1}$	Kinetic parameter for AOB activity	0.025 (0.024)	h
$\theta_{AOB,2}$	Kinetic parameter for AOB activity	0.1 (0.13)	$\text{h}\cdot\text{L}\cdot\text{mol}^{-1}$
$\theta_{AOB,3}$	Kinetic parameter for AOB activity	2.5 (2.4178)	$\text{h}\cdot\text{L}^2\cdot\text{mol}^{-2}$
$\theta_{NOB,1}$	Kinetic parameter for NOB activity	$0.11 \cdot 10^{-3}$ ($0.1 \cdot 10^{-3}$)	h
$\theta_{NOB,2}$	Kinetic parameter for NOB activity	1.1 (1.13)	$\text{h}\cdot\text{L}\cdot\text{mol}^{-1}$
$\Lambda (\Lambda_h)$	Extent variance-covariance matrix (for the h th sample)	–	mol^2
$\lambda (\lambda_{i,h})$	Extent variance (for the i th reaction and the h th sample)	–	mol^2
$\Sigma (\Sigma_h)$	Measurement error variance-covariance matrix (for the h th sample)	–	
σ_{TAN}	Measurement standard deviation for TAN	0.01	$\text{mol} \cdot \text{L}^{-1}$
σ_{TNO2}	Measurement standard deviation for TNO2	0.01	$\text{mol} \cdot \text{L}^{-1}$
σ_{TNO3}	Measurement standard deviation for TNO3	0.01	$\text{mol} \cdot \text{L}^{-1}$
σ_{pH}	Measurement standard deviation for pH	0.05	–
τ	Integrand (time)	–	h
$[\cdot]$	concentration symbol equivalent to c	–	$\text{mol} \cdot \text{L}^{-1}$

Table 2: List of candidate rate laws used for both nitritation and nitrataion reactions. The substrate concentration c_S is the free ammonia concentration $[\text{NH}_3]$ for the nitritation and the free nitrous acid concentration $[\text{HNO}_2]$ for the nitrataion.

Name	Index	Candidate rate law	Parameter vector
	j	$r_{k,i}^{(j)}(c_S, \boldsymbol{\theta}_i^{(j)}), i \in \{1, 2\}$	$\boldsymbol{\theta}_i^{(j)}$
Zeroth order	1	$\begin{cases} 1/\theta_{i,1}^{(1)} & \text{if } c_S \geq 0 \\ 0 & \text{otherwise} \end{cases}$	$[\theta_{i,1}^{(1)}]$
First order	2	$\frac{c_S}{\theta_{i,1}^{(2)}}$	$[\theta_{i,1}^{(2)}]$
Monod	3	$\frac{c_S}{\theta_{i,1}^{(3)} + \theta_{i,2}^{(3)} c_S}$	$[\theta_{i,1}^{(3)} \quad \theta_{i,2}^{(3)}]^T$
Tessier	4	$\frac{1 - \exp(-c_S \theta_{i,2}^{(4)} / \theta_{i,1}^{(4)})}{\theta_{i,2}^{(4)}}$	$[\theta_{i,1}^{(4)} \quad \theta_{i,2}^{(4)}]^T$
Haldane	5	$\frac{c_S}{\theta_{i,1}^{(5)} + \theta_{i,2}^{(5)} c_S + \theta_{i,3}^{(5)} c_S^2}$	$[\theta_{i,1}^{(5)} \quad \theta_{i,2}^{(5)} \quad \theta_{i,3}^{(5)}]^T$

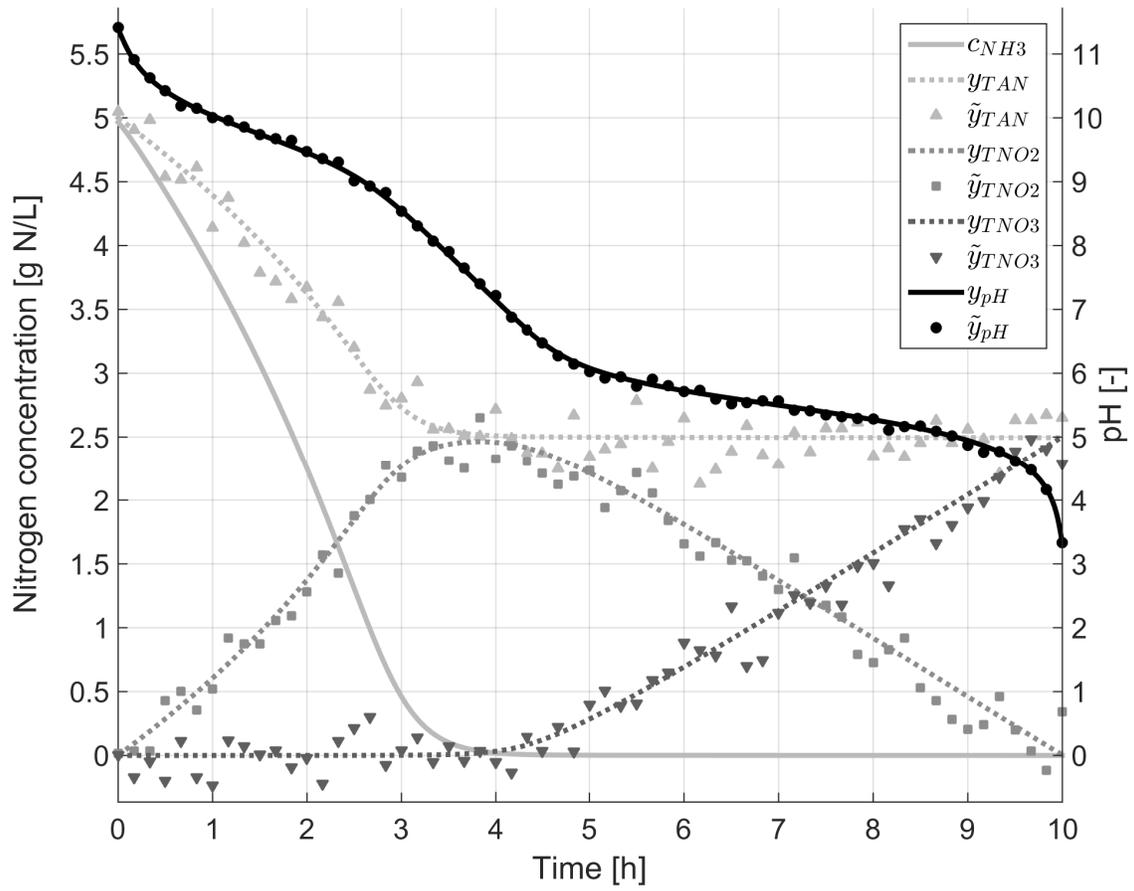


Figure 1: **Data generation:** Simulated concentration (in gN/L) and pH (continuous and dashed lines) with corresponding measurements (dots, squares, triangles).

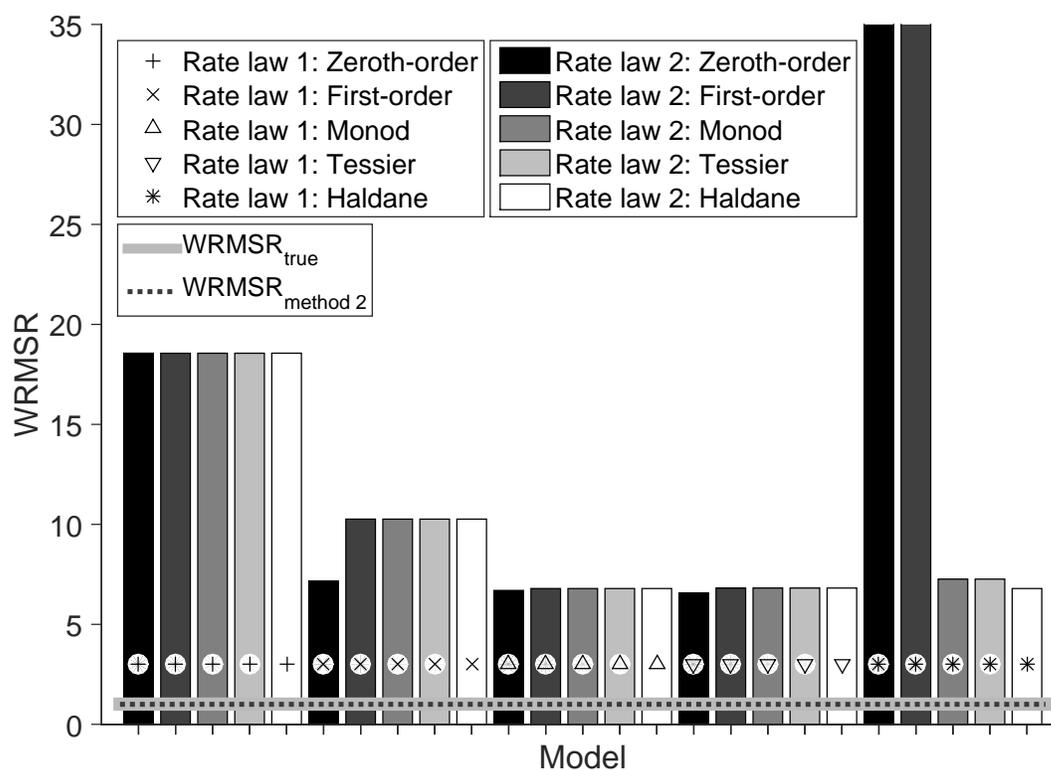


Figure 2: **Method 1: simultaneous model identification.** WRMSR values for 25 models. The markers indicate the selected rate law for the first rate law. Shading of the bars indicates the selected rate law for the second reaction. The WRMSR values for the true model and for the best model obtained with Method 2 are indicated by a full and a dashed line, respectively.

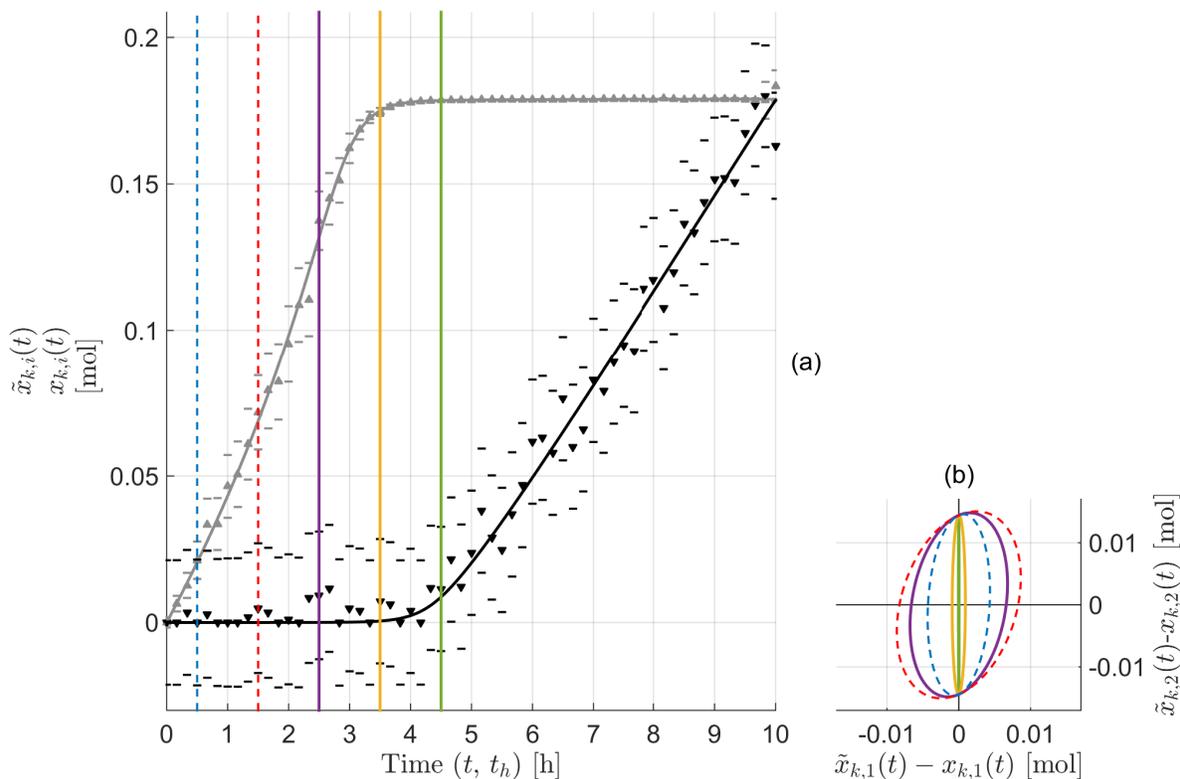


Figure 3: **Method 2: extent-based modeling – Step 1: Computation of experimental extents.** (a) True (lines) and experimental (dots) extents with 3σ confidence intervals for the nitritation (gray) and nitration (black) reactions. (b) Variance-covariance matrix as 3σ confidence region (ellipsoid) for the experimental extent errors around (0,0); colored lines corresponding to $t_h = 0.5, 1.5, 2.5, 3.5$ and 4.5 h indicated with matching colors and styles in (a).

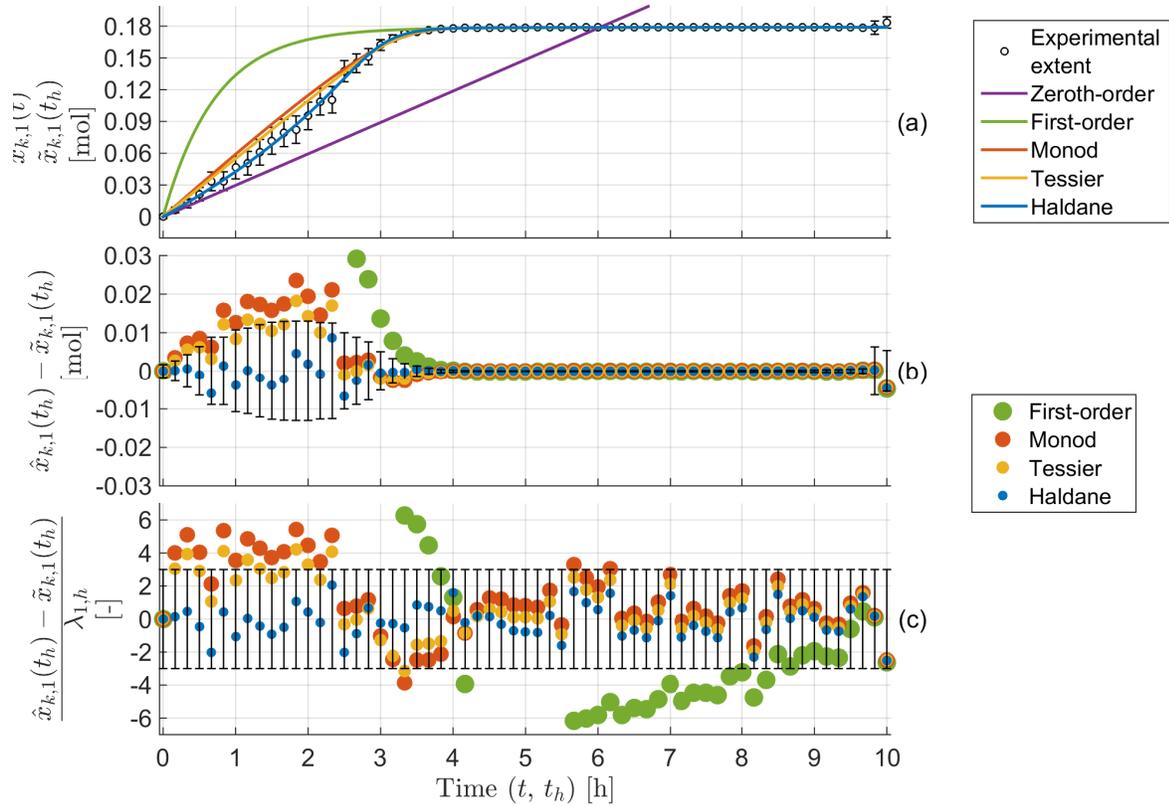


Figure 4: **Method 2: Extent-based modeling – Step 2: Modeling of the 1st extent.** (a) Experimental (circles, with error bars) and modeled (continuous lines) extents as functions of time; (b) Residuals between modeled and experimental extents as functions of time; (c) Normalized residuals between modeled and experimental extents as functions of time.

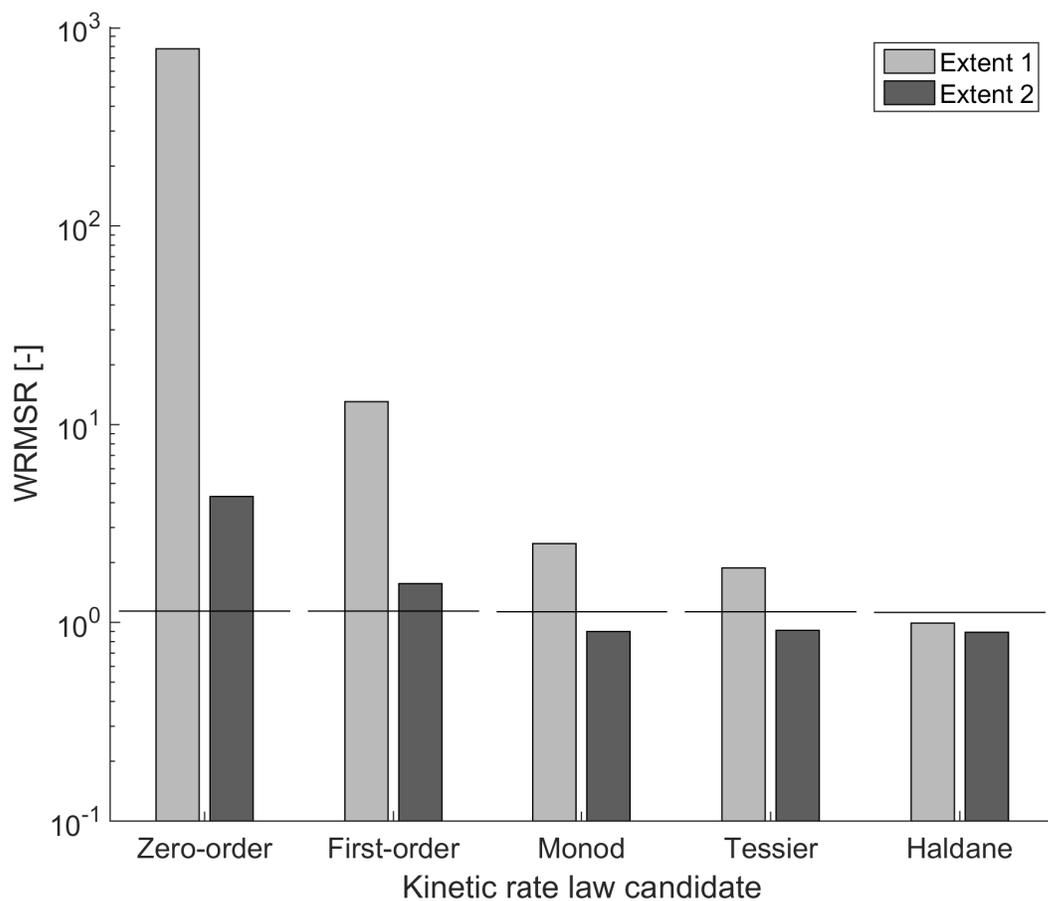


Figure 5: **Method 2: extent-based modeling – Step 2: Modeling of extents – Lack-of-Fit.** WRMSR for all extents and all candidate rate laws (bars) and 95% upper control limits of the associated χ^2 -distribution (lines).

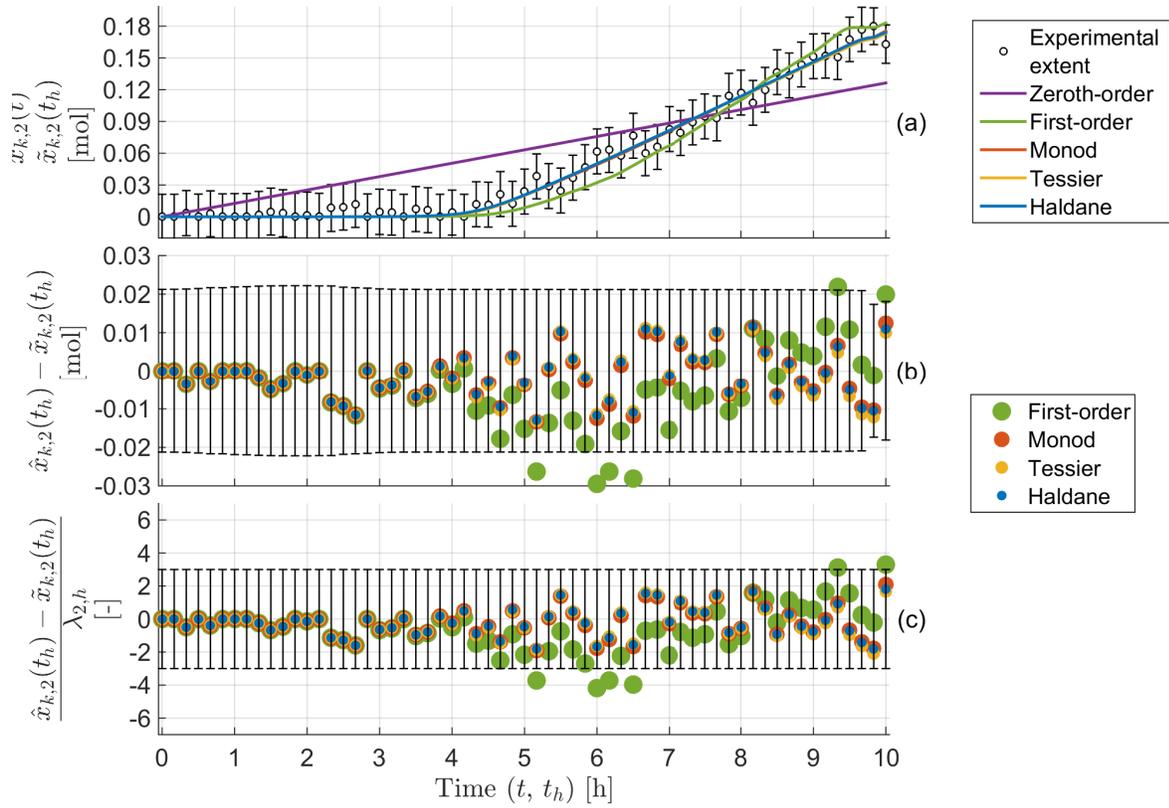


Figure 6: **Method 2: extent-based modeling – Step 2: Modeling of the 2nd extent.** (a) Experimental (circles, with error bars) and modeled (continuous lines) extents as functions of time; (b) Residuals between modeled and experimental extents as functions of time; (c) Normalized residuals between modeled and experimental extents as functions of time.

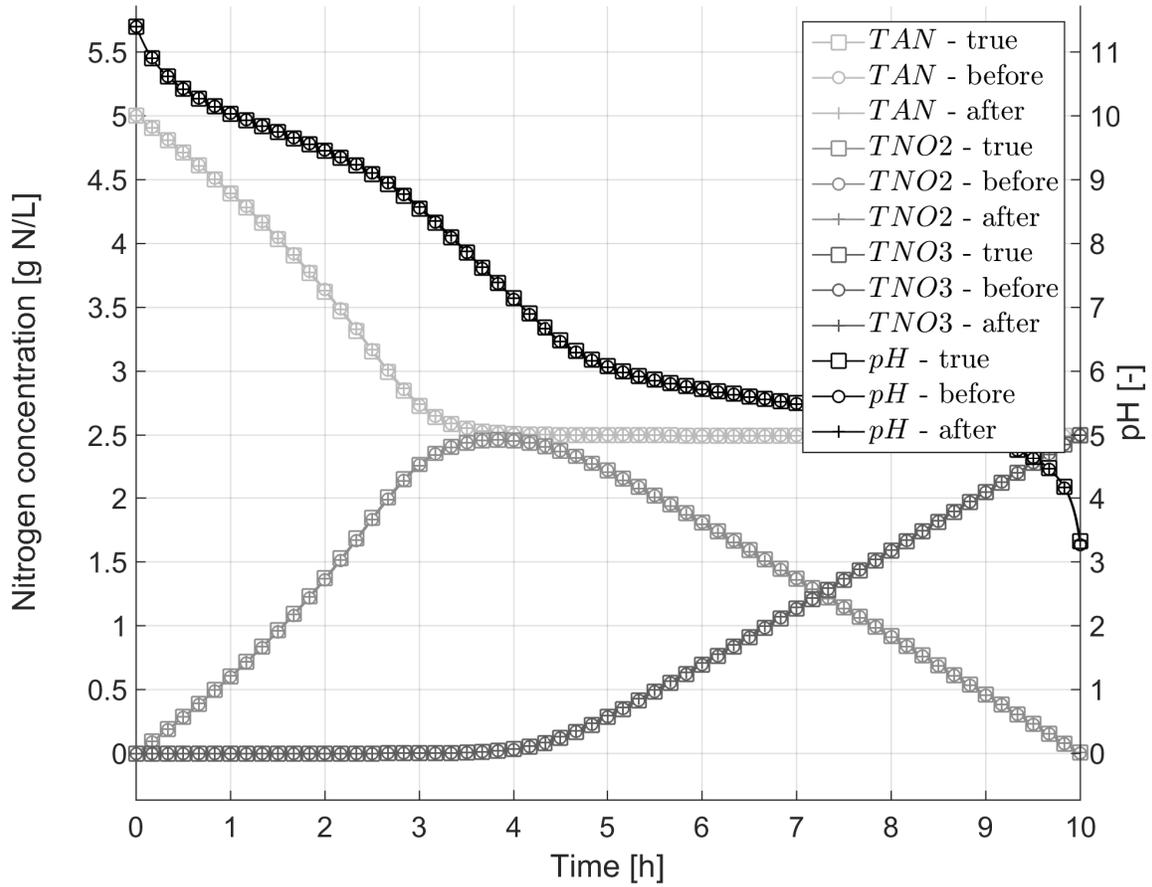


Figure 7: **Method 2: extent-based modeling – Step 3: Model fine-tuning.** Simulation of the TAN, TNO2, and TNO3 concentrations and pH for (i) the true data-generating model, (ii) the model obtained before fine-tuning, and (iii) the model obtained after fine-tuning. Differences between these simulations are barely noticeable.

Supporting Information Available

Supporting Information includes bounding procedures, additional figures, and all code to produce our results.

This material is available free of charge via the Internet at <http://pubs.acs.org/>.

References

- (1) Ni, B. J.; Peng, L.; Law, Y.; Guo, J.; Yuan, Z. Modeling of nitrous oxide production by autotrophic ammonia-oxidizing bacteria with multiple production pathways. *Environmental Science & Technology* **2014**, *48*, 3916–3924.
- (2) Liu, L.; Binning, P. J.; Smets, B. F. Evaluating alternate biokinetic models for trace pollutant cometabolism. *Environmental Science & Technology* **2015**, *49*, 2230–2236.
- (3) Henze, M.; Gujer, W.; Mino, T.; van Loosdrecht, M. *Activated sludge models ASM1, ASM2, ASM2d and ASM3. IWA Scientific and Technical Report.*; IWA Publishing, London, UK, 2000.
- (4) Jenkins, D.; Wanner, J. Activated Sludge – 100 Years and Counting. *Water Intelligence Online* **2014**, *13*.
- (5) Brun, R.; Reichert, P.; Künsch, H. R. Practical identifiability analysis of large environmental simulation models. *Water Resources Research* **2001**, *37*, 1015–1030.
- (6) Liu, C.; Zachara, J. M. Uncertainties of Monod kinetic parameters nonlinearly estimated from batch experiments. *Environmental Science & Technology* **2001**, *35*, 133–141.
- (7) Nelly, N.; Müller, T. G.; Gyllenberg, M.; Timmer, J. Quantitative analyses of anaerobic wastewater treatment processes: identifiability and parameter estimation. *Biotechnology and Bioengineering* **2002**, *78*, 89–103.

- 490 (8) Brockmann, D.; Rosenwinkel, K.-H.; Morgenroth, E. Practical identifiability of bioki-
491 netic parameters of a model describing two-step nitrification in biofilms. *Biotechnology*
492 *and Bioengineering* **2008**, *101*, 497–514.
- 493 (9) Neumann, M. B.; Gujer, W. Underestimation of uncertainty in statistical regression of
494 environmental models: influence of model structure uncertainty. *Environmental Science*
495 *& Technology* **2008**, *42*, 4037–4043.
- 496 (10) Bennett, N. D.; Croke, B. F.; Guariso, G.; Guillaume, J. H.; Hamilton, S. H.; Jake-
497 man, A. J.; Marsili-Libelli, S.; Newham, L. T.; Norton, J. P.; Perrin, C.; Pierce, S. A.
498 Characterising performance of environmental models. *Environmental Modelling & Soft-*
499 *ware*, **2013**, *40*, 1–20.
- 500 (11) Brugnach, M.; Pahl-Wostl, C.; Lindenschmidt, K.; Janssen, J.; Filatova, T.; Mou-
501 ton, A.; Holtz, G.; Van der Keur, P.; Gaber, N. *Developments in Integrated Environ-*
502 *mental Assessment, Vol.3*; Elsevier, 2008; Chapter 4 - Complexity and uncertainty:
503 Rethinking the modelling activity, pp 49–68.
- 504 (12) Sin, G.; Ödman, P.; Petersen, N.; Lantz, A. E.; Gernaey, K. V. Matrix notation for
505 efficient development of first-principles models within PAT applications: Integrated
506 modeling of antibiotic production with *Streptomyces coelicolor*. *Biotechnology and Bio-*
507 *engineering* **2008**, *101*, 153–171.
- 508 (13) Jakeman, A. J.; Letcher, R. A.; Norton, J. P. Ten iterative steps in development and
509 evaluation of environmental models. *Environmental Modelling & Software* **2006**, *21*,
510 602–614.
- 511 (14) Müller, T.; Dürr, R.; Isken, B.; Schulze-Horsel, J.; Reichl, U.; Kienle, A. Distributed
512 modeling of human influenza a virus-host cell interactions during vaccine production.
513 *Biotechnology and Bioengineering* **2013**, *110*, 2252–2266.

- 514 (15) Zambrano-Bigiarini, M.; Rojas, R. A model-independent Particle Swarm Optimisation
515 software for model calibration. *Environmental Modelling & Software* **2013**, *43*, 5–25.
- 516 (16) Larsen, T. A., Udert, K. M., Lienert, J., Eds. *Source separation and decentralization*
517 *for wastewater management*; IWA Publishing, 2013.
- 518 (17) Udert, K. M.; Wächter, M. Complete nutrient recovery from source-separated urine by
519 nitrification and distillation. *Water research* **2012**, *46*, 453–464.
- 520 (18) Rieger, L.; Gillot, S.; Langergraber, G.; Ohtsuki, T.; Shaw, A.; Takács, I.; Winkler, S.
521 *Guidelines for using activated sludge models. IWA Task Group on Good Modelling Prac-*
522 *tice. IWA Scientific and Technical Report*; IWA Publishing., 2012.
- 523 (19) Van De Steene, M.; Van Vooren, L.; Ottoy, J. P.; Vanrolleghem, P. A. Automatic buffer
524 capacity model building for advanced interpretation of titration curves. *Environmental*
525 *Science & Technology* **2002**, *36*, 715–723.
- 526 (20) Mašić, A.; Udert, K.; Villez, K. Global parameter optimization for biokinetic modeling
527 of simple batch experiments. *Environmental Modelling and Software* **2016**, *85*, 356–373.
- 528 (21) Bhatt, N.; Amrhein, M.; Bonvin, D. Incremental Identification of Reaction and Mass-
529 Transfer Kinetics Using the Concept of Extents. *Industrial & Engineering Chemistry*
530 *Research* **2011**, *50*, 12960–12974.
- 531 (22) Srinivasan, S.; Billeter, J.; D., B. Extent-based incremental identification of reaction
532 systems using concentration and calorimetric measurements. *Chemical Engineering*
533 *Journal* **2012**, *207-208*, 785–793.
- 534 (23) Billeter, J.; Srinivasan, S.; D., B. Extent-based Kinetic Identification using Spectro-
535 scopic Measurements and Multivariate Calibration. *Analytica Chimica Acta* **2013**, *767*,
536 21–34.

- 537 (24) Srinivasan, S.; Billeter, J.; D., B. Sequential Model Identification of Reaction Systems
538 - The Missing Path between the Incremental and Simultaneous Approaches. *AIChE*
539 *Journal* **2017**, submitted.
- 540 (25) Fumasoli, A. Nitrification of Urine as Pretreatment for Nutrient Recovery. Ph.D. thesis,
541 ETH Zürich, 2016.
- 542 (26) Srinivasan, S.; Billeter, J.; Bonvin, D. Identification of Multiphase Reaction Systems
543 with Instantaneous Equilibria. *Industrial & Engineering Chemistry Research* **2016**, *29*,
544 8034–8045.
- 545 (27) Alonso, C.; Zhu, X.; Suidan, M. T.; Kim, B. R.; Kim, B. J. Parameter estimation in
546 biofilter systems. *Environmental Science & Technology* **2000**, *34*, 2318–2323.
- 547 (28) Zhou, Y. A. N.; Pijuan, M.; Zeng, R. J.; Yuan, Z. Free nitrous acid inhibition on ni-
548 trous oxide reduction by a denitrifying-enhanced biological phosphorus removal sludge.
549 *Environmental Science & Technology* **2008**, *42*, 8260–8265.
- 550 (29) Ni, B. J.; Zeng, R. J.; Fang, F.; Xu, J.; Sheng, G. P.; Yu, H. Q. A novel approach
551 to evaluate the production kinetics of extracellular polymeric substances (EPS) by
552 activated sludge using weighted nonlinear least-squares analysis. *Environmental Science*
553 *& Technology* **2009**, *43*, 3743–3750.
- 554 (30) Sathyamoorthy, S.; Chandran, K.; Ramsburg, C. A. Biodegradation and cometabolic
555 modeling of selected beta blockers during ammonia oxidation. *Environmental Science*
556 *& Technology* **2013**, *47*, 12835–12843.
- 557 (31) Nelder, J. A.; Mead, R. A simplex method for function minimization. *The Computer*
558 *Journal* **1965**, *7*, 308–313.
- 559 (32) Chandran, K.; Smets, B. F. Estimating biomass yield coefficients for autotrophic am-

- 560 monia and nitrite oxidation from batch respirograms. *Water Research* **2001**, *35*, 3153–
561 3156.
- 562 (33) Buendía, I. M.; Fernández, F. J.; Villaseñor, J.; Rodríguez, L. Feasibility of anaerobic
563 co-digestion as a treatment option of meat industry wastes. *Bioresource Technology*
564 **2009**, *100*, 1903–1909.
- 565 (34) Gernaey, K.; Bogaert, H.; Massone, A.; Vanrolleghem, P.; Verstraete, W. On-line nitrifi-
566 cation monitoring in activated sludge with a titrimetric sensor. *Environmental Science*
567 *& Technology* **1997**, *31*, 2350–2355.
- 568 (35) Jensen, P. D.; Ge, H.; Batstone, D. J. Assessing the role of biochemical methane po-
569 tential tests in determining anaerobic degradability rate and extent. *Water Science &*
570 *Technology* **2011**, *64*, 880–886.
- 571 (36) Santa Cruz, J. A.; Mussati, S. F.; Scenna, N. J.; Gernaey, K. V.; Mussati, M. C.
572 Reaction invariant-based reduction of the activated sludge model ASM1 for batch ap-
573 plications. *Journal of Environmental Chemical Engineering* **2016**, *4*, 3654–3664.
- 574 (37) Rodrigues, D.; Srinivasan, S.; Billeter, J.; D., B. Variant and Invariant States for Chem-
575 ical Reaction Systems. *Computers & Chemical Engineering* **2015**, *73*, 23–33.
- 576 (38) Bonvin, D.; Rippin, D. W. T. Target factor analysis for the identification of stoichio-
577 metric models. *Chemical Engineering Science* **1990**, *45*, 3417–3426.
- 578 (39) Mašić, A.; ; Billeter, J.; Bonvin, D.; Villez, K. Extent computation under rank-deficient
579 conditions. *20th World Congress of the International Federation of Automatic Control*
580 *(IFAC2017), 9-14 July 2017, Toulouse, France* **2017**, Accepted for oral presentation.
- 581 (40) Rieger, L.; Alex, J.; Winkler, S.; Boehler, M.; Thomann, M.; Siegrist, H. Progress in
582 sensor technology – progress in process control? Part I: Sensor property investigation
583 and classification. *Water Science & Technology* **2003**, *47*, 103–112.

- 584 (41) Rosén, C.; Rieger, L.; Jeppsson, U.; Vanrolleghem, P. A. Adding realism to simulated
585 sensors and actuators. *Water Science & Technology* **2008**, *57*, 337–344.
- 586 (42) Schielke-Jenni, S.; Villez, K.; Morgenroth, E.; Udert, K. M. Observability of anammox
587 activity in single-stage nitrification/anammox reactors using mass balances. *Environ-
588 mental Science: Water Research & Technology* **2015**, *1*, 523–534.
- 589 (43) Mašić, A.; Srinivasan, S.; Billeter, J.; Bonvin, D.; Villez, K. Shape Constrained Splines
590 as Transparent Black-Box Models for Bioprocess Modeling. *Computers and Chemical
591 Engineering* **2017**, *99*, 96–105.
- 592 (44) The MathWorks Inc., *MATLAB Release 2014b*; Natick, Massachusetts, 2014.
- 593 (45) Villez, K.; Rengaswamy, R.; Venkatasubramanian, V. Generalized Shape Constrained
594 Spline Fitting for Qualitative Analysis of Trends. *Computers & Chemical Engineering*
595 **2013**, *58*, 116–134.

596 **Summary**

597 The Supplementary Information consists of:

- 598 • This text which consists of 34 pages and includes 27 figures.
- 599 • The latest version of the EMI software package which enables reproduction of our
600 results in the Matlab environment.

601 **Software**

602 All software necessary to reproduce the results presented in this work is available as part of
603 the self-sufficient Efficient Model Identification (EMI) package for Matlab or Octave. It is
604 published under the GPL v3 open-source license and constitutes the Supporting Information
605 together with this text.

606 **Graphical overview of modeling via extents**

607 The modeling procedure is illustrated in Fig. S.1 for the exemplary case studied in this
608 work. The three main steps, i.e. *(i)* extent computation, *(ii)* extent modeling, and *(iii)*
609 model fine-tuning, are shown from top to bottom. The experimental extents are split into
610 two individual time series corresponding to the two reactions. After this, the parameters for
611 four candidate rate laws are estimated for each reaction. The best-fit rate laws are combined
612 into a joint model. The associated parameter values are used as an initial guess for the
613 fine-tuning step.

STEP 1 – EXTENT COMPUTATION

\tilde{y}
Measurements

\tilde{x}_k
Experimental extents

STEP 2 – EXTENT MODELING

$\tilde{x}_{k,1}$
Experimental extent 1

$\tilde{x}_{k,2}$
Experimental extent 2

$\hat{\theta}_{k,1}^{(1)}$ Rate law 1 $\hat{\theta}_{k,1}^{(2)}$ Rate law 2 $\hat{\theta}_{k,1}^{(3)}$ Rate law 3 $\hat{\theta}_{k,1}^{(4)}$ Rate law 4 $\hat{\theta}_{k,1}^{(5)}$ Rate law 5 $\hat{\theta}_{k,2}^{(1)}$ Rate law 1 $\hat{\theta}_{k,2}^{(2)}$ Rate law 2 $\hat{\theta}_{k,2}^{(3)}$ Rate law 3 $\hat{\theta}_{k,2}^{(4)}$ Rate law 4 $\hat{\theta}_{k,2}^{(5)}$ Rate law 5

$\hat{r}_{k,1}$
Best-fit rate law

$\hat{r}_{k,2}$
Best-fit rate law

STEP 3 – MODEL FINE-TUNING

\hat{r}_k
Initial guess

\hat{r}_k
Final model

Figure S.1: **Extent-based modeling procedure.** Through the computation of experimental extents, kinetic modeling can be divided in smaller problems, each one focusing on the identification of the rate law and the corresponding parameters for a single reaction. A fine-tuning step is used at the end to obtain the final parameter estimates for the identified rate laws.

614 Bounding procedures

615 The estimation of the kinetic parameters in step 2 of the incremental model identification
616 procedure is based on the branch-and-bound algorithm. Its use and application for bioki-
617 netic model parameter estimation has been demonstrated before²⁰. The following bounding
618 procedures constitute the only differences with this prior work. In what follows, we consider
619 the estimation of a single parameter vector, $\boldsymbol{\theta}^{(j)}$, for a single candidate reaction rate law,
620 $r_{k,i}^{(j)}$. For the sake of conciseness, these are given as $\boldsymbol{\theta}$ and r in what follows.

621 Definition of considered parameter set and parameter subsets

During the branch-and-bound algorithm, several hyper-rectangular parameter subsets are considered. These subsets are denoted here as Ω_a with a an integer indicating the chronology of the evaluated parameter subsets. Ω_0 corresponds to the root set, i.e. the set containing all feasible parameter values. Each parameter subset can be described as follows:

$$\boldsymbol{\theta} \in \Omega_a \Leftrightarrow \boldsymbol{\theta}_a^L \leq \boldsymbol{\theta} \leq \boldsymbol{\theta}_a^U \quad (38)$$

622 with $\boldsymbol{\theta}_a^L$ and $\boldsymbol{\theta}_a^U$ containing the lower and upper bounds for each element of $\boldsymbol{\theta}$. Inequalities
623 between vectors are defined in an element-wise manner.

624 Upper bound

625 An upper bound to the objective function value is easily obtained by evaluating the objective
626 function in (32) at an arbitrary value for $\boldsymbol{\theta}$ within the considered set Ω_a :

$$Q^U = q(\boldsymbol{\theta}) = \sum_{h=1}^H \frac{(\tilde{x}_{k,i}(t_h) - x_{k,i}(t_h))^2}{\lambda_{i,h}}, \quad (39)$$

627 following simulation of the following DAE system:

$$\mathbf{g}(\mathbf{n}(t)/V) = \mathbf{0} \quad (40)$$

$$\bar{\mathbf{E}} \mathbf{n}(t) = \bar{\mathbf{n}}_{k,0} + \bar{\mathbf{N}}^T \mathbf{x}_k(t) \quad (41)$$

$$\forall r = 1, \dots, R_k : x_{k,r}(t) = \begin{cases} V \int_0^t r_{k,i}(\mathbf{n}(\tau)/V, \boldsymbol{\theta}) d\tau, & x_{k,i}(0) = 0 \quad \text{if } r = i \\ \mathcal{I}(\mathbf{t}, \tilde{\mathbf{x}}_{k,r}, t), & \text{if } r \neq i \end{cases} \quad (42)$$

628 with definitions as in the main text. It is fairly trivial to see that Q^U is a valid upper
 629 bound. Indeed, at least one set of parameter values within Ω_a results in an objective value
 630 that is lower or equal to Q^U . This is true since the evaluated $\boldsymbol{\theta}$ is in the set and gives
 631 $q(\boldsymbol{\theta}) = Q^U$.

632 Lower bound

633 As usual, obtaining a provable lower bound is more challenging. In this study, we follow the
 634 previously developed procedure.²⁰ The main difference is that there is no need *(i)* to apply
 635 a linearizing model reformulation or *(ii)* to bound the value of rate measurements. This is
 636 because *(i)* the initial conditions and the stoichiometric matrix are considered known at the
 637 stage of kinetic parameter estimation and *(ii)* extents are integral states. As a result, the
 638 bounding procedure remains fairly simple.

639 To start, consider that the reaction rate can be bounded as follows for each of the can-
 640 didate rate laws (see main text, Table 2):

$$\boldsymbol{\theta} \in \Omega_a : 0 \leq r(\mathbf{n}(\tau)/V, \boldsymbol{\theta}_a^U) \leq r(\mathbf{n}(\tau)/V, \boldsymbol{\theta}) \leq r(\mathbf{n}(\tau)/V, \boldsymbol{\theta}_a^L) \quad (43)$$

641 Indeed, thanks to the particular parameterization in Table 2, one can easily see that the
 642 highest (lowest) reaction rates are obtained for the lowest (highest) parameter values within

643 Ω_a , i.e. $\boldsymbol{\theta}_a^L$ ($\boldsymbol{\theta}_a^U$). In addition, the reaction rate is strictly non-negative at all times (i.e.
 644 irreversible reaction). Combining this positivity of the reaction rate with the bounds for the
 645 reaction rates means that one can write the following inequality for i th modeled extent of
 646 reaction:

$$\boldsymbol{\theta} \in \Omega_a : \quad x_{k,i}^L(t) \leq x_{k,i}(t) \leq x_{k,i}^U(t) \quad (44)$$

647 with

$$x_{k,i}^L(t) = V \int_0^t r(\mathbf{n}(\tau)/V, \boldsymbol{\theta}_a^U) d\tau, \quad x_{k,i}^L(0) = 0 \quad (45)$$

$$x_{k,i}^U(t) = V \int_0^t r(\mathbf{n}(\tau)/V, \boldsymbol{\theta}_a^L) d\tau, \quad x_{k,i}^U(0) = 0 \quad (46)$$

648 subject to (40–41) and all evaluations of the $r \neq i$ case in (42).

649 In words, the considered extent of reaction at time t is the highest (lowest) for the highest
 650 (lowest) reaction rates and thus the lowest (highest) parameter values. This statement follows
 651 from the fact that the extent is a monotonic function of time (positivity of the reaction rate)
 652 with the derivative defined by the reaction rate. This derivative takes its lowest (highest)
 653 attainable value for the highest (lowest) parameter values at any time t and for any possible
 654 state that has been reached at time t . It follows that two simulations delivering $x_{k,i}^U(t)$ and
 655 $x_{k,i}^L(t)$ deliver effective bounds to the extent profiles obtained with any feasible value for $\boldsymbol{\theta}$
 656 within Ω_a .

657 Based on interval arithmetic, the squared residuals $s_{k,i}(t_h) = (\tilde{x}_{k,i}(t_h) - x_{k,i}(t_h))^2$ can
 658 now be lower bounded as follows:

$$\boldsymbol{\theta} \in \Omega_a : s_{k,i}^L(t_h) \leq s_{k,i}(t_h)|_{\boldsymbol{\theta}} \quad (47)$$

659 with

$$d_{k,i}^L(t_h) = \tilde{x}_{k,i}(t_h) - x_{k,i}^U(t_h) \quad (48)$$

$$d_{k,i}^U(t_h) = \tilde{x}_{k,i}(t_h) - x_{k,i}^L(t_h) \quad (49)$$

$$s_{k,i}^L(t_h) = \begin{cases} 0 & \text{if } d_{k,i}^L(t_h) \leq 0 \leq d_{k,i}^U(t_h) \\ \min \left(d_{k,i}^L(t_h)^2, d_{k,i}^U(t_h)^2 \right) & \text{otherwise} \end{cases} \quad (50)$$

660 From this, it follows that:

$$\boldsymbol{\theta} \in \Omega_a : Q^L \leq q(\boldsymbol{\theta}) \quad (51)$$

661 with

$$Q^L = \sum_{h=1}^H \frac{s_{k,i}^L(t_h)}{\lambda_{i,h}}, \quad (52)$$

662 which proves that Q^L is a valid lower bound.

663 **Implementation of the bounding procedures**

664 The above procedures suggest simulation of the considered extent of reaction for three pa-
 665 rameter vectors. The first is executed for an arbitrary feasible choice for $\boldsymbol{\theta}$ within Ω_a . The
 666 second and third simulation is executed for $\boldsymbol{\theta}^L$ and $\boldsymbol{\theta}^U$. These simulations are the computa-

tionally most expensive steps of the bounding procedures. For this reason, the upper bound procedure is evaluated for θ^L and θ^U , since the corresponding extent simulations are required anyway for the lower bound. This means only two simulations are executed to compute both the lower and upper bound. In the process, one obtains two distinct upper bound values. The minimum of these two upper bounds is then reported as the best-known upper bound. A graphical scheme of the bounding procedures is given in Fig. S.2. Note that this scheme is fairly simple compared to the original bounding procedures.²⁰

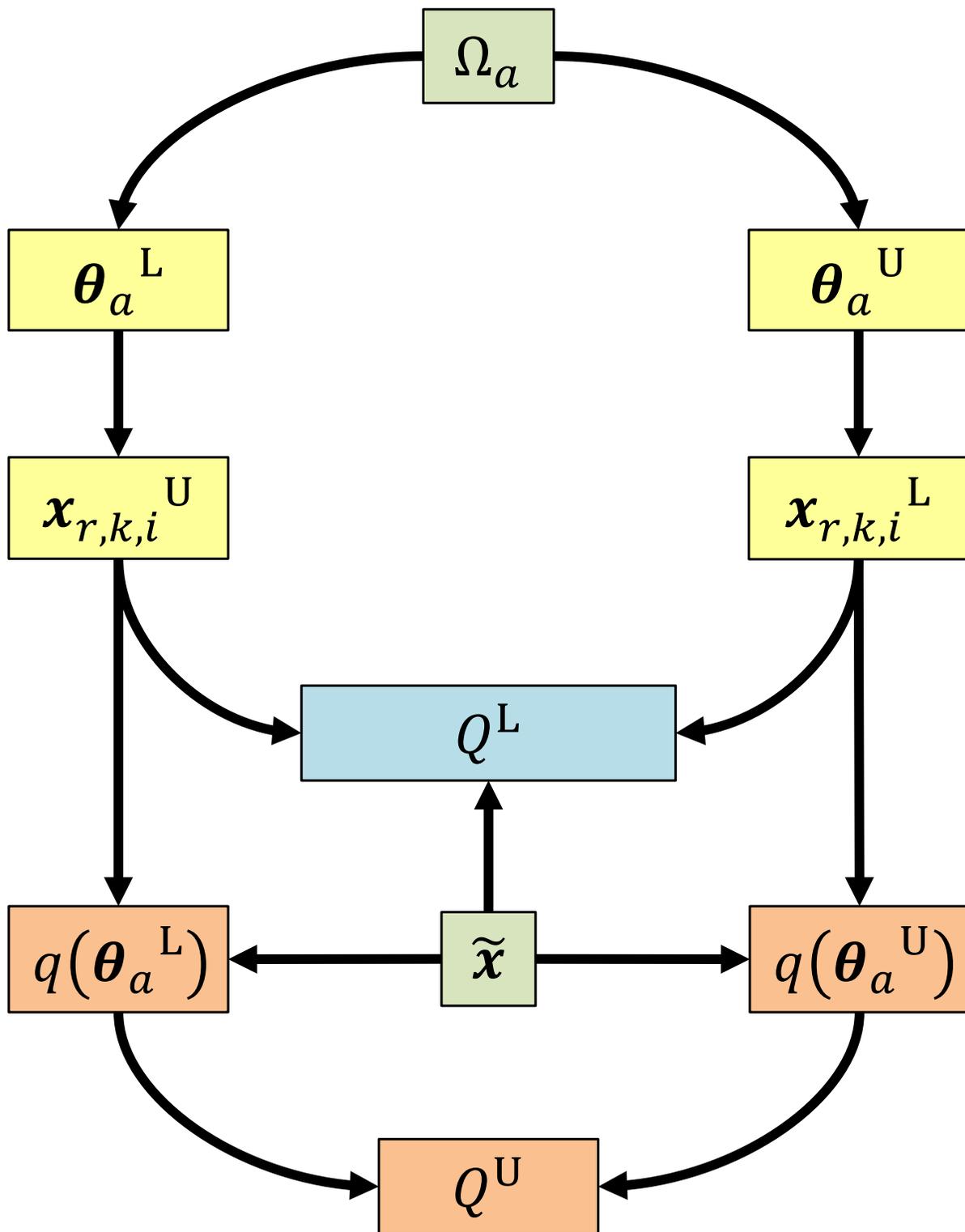


Figure S.2: **Illustration of the bounding procedures:** Two simulations are executed, one for both extremal parameter vectors (θ_a^L , θ_a^U) that bound the considered set (Ω_a). These deliver the bounding profiles for the extent of reaction ($x_{k,i}^L$ and $x_{k,i}^U$). By combining these two extremal profiles with the experimental extent series (\tilde{x}), one can compute both the upper bound (Q^U) and the lower bound (Q^L).
 In this scheme, should $Q(\theta_a^L)$ and $Q(\theta_a^U)$ not be $q(\theta_a^L)$ and $Q(\theta_a^U)$? Also, the exponents L and U seem to be bold although they should not be...

674 **Additional results**

675 The next figures (Fig. S.3-S.27) show the simulation results obtained with each model
676 obtained with Method 1 after parameter estimation.

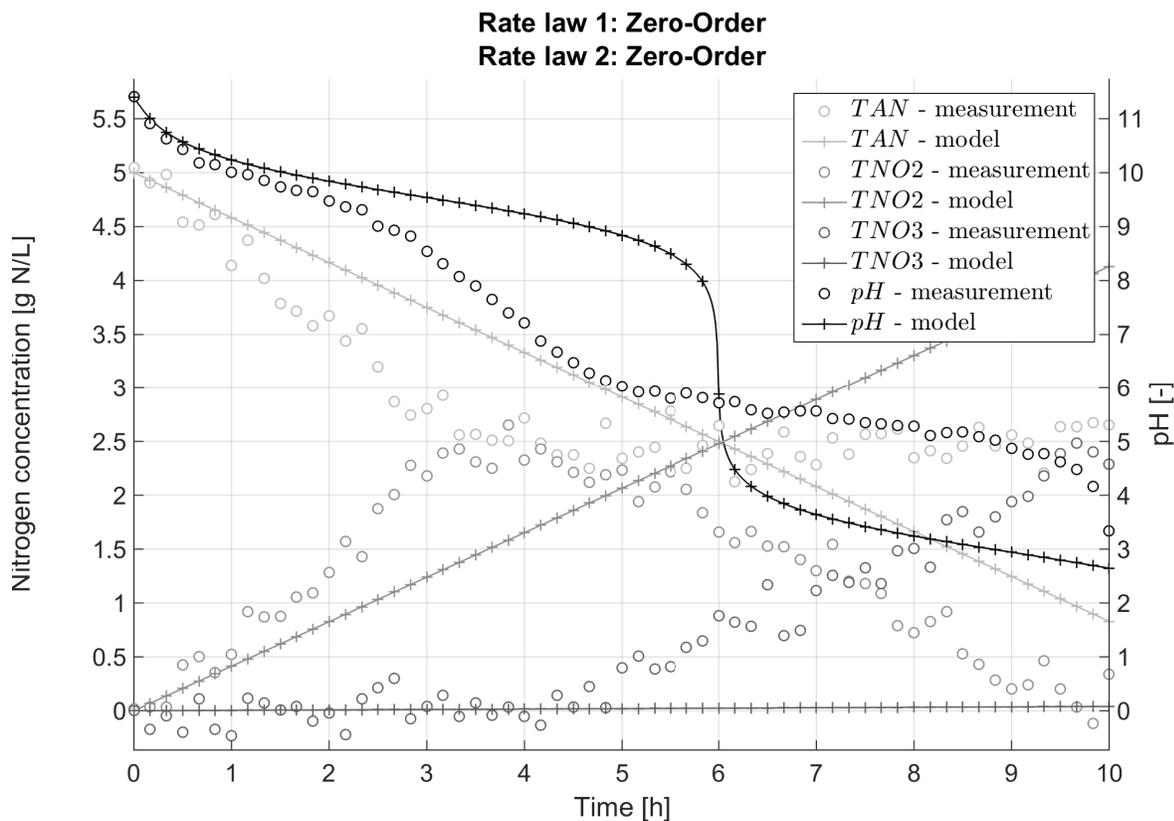


Figure S.3: **Method 1 - Simultaneous model identification - Model 1.** Measurements and simulation of the measured variables with model 1 after parameter estimation with the Nelder-Mead simplex method.

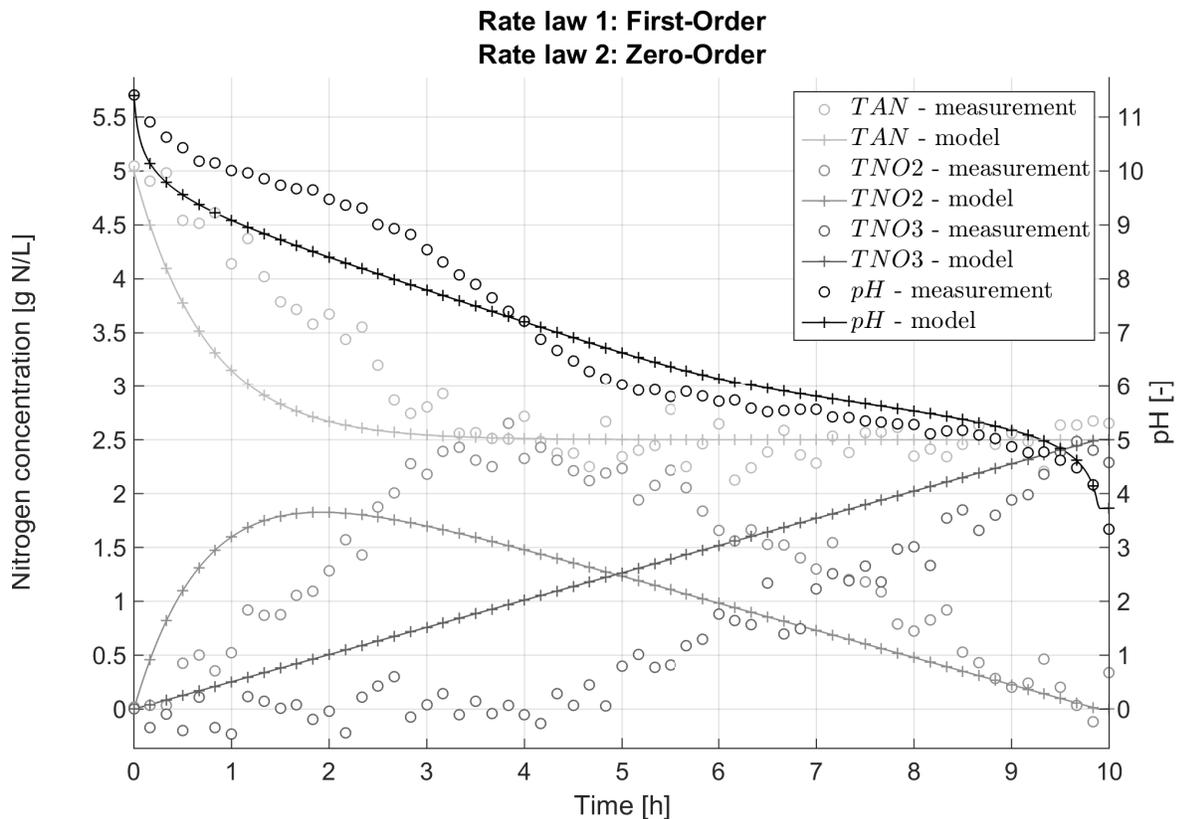


Figure S.4: **Method 1 - Simultaneous model identification - Model 2.** Measurements and simulation of the measured variables with model 2 after parameter estimation with the Nelder-Mead simplex method.

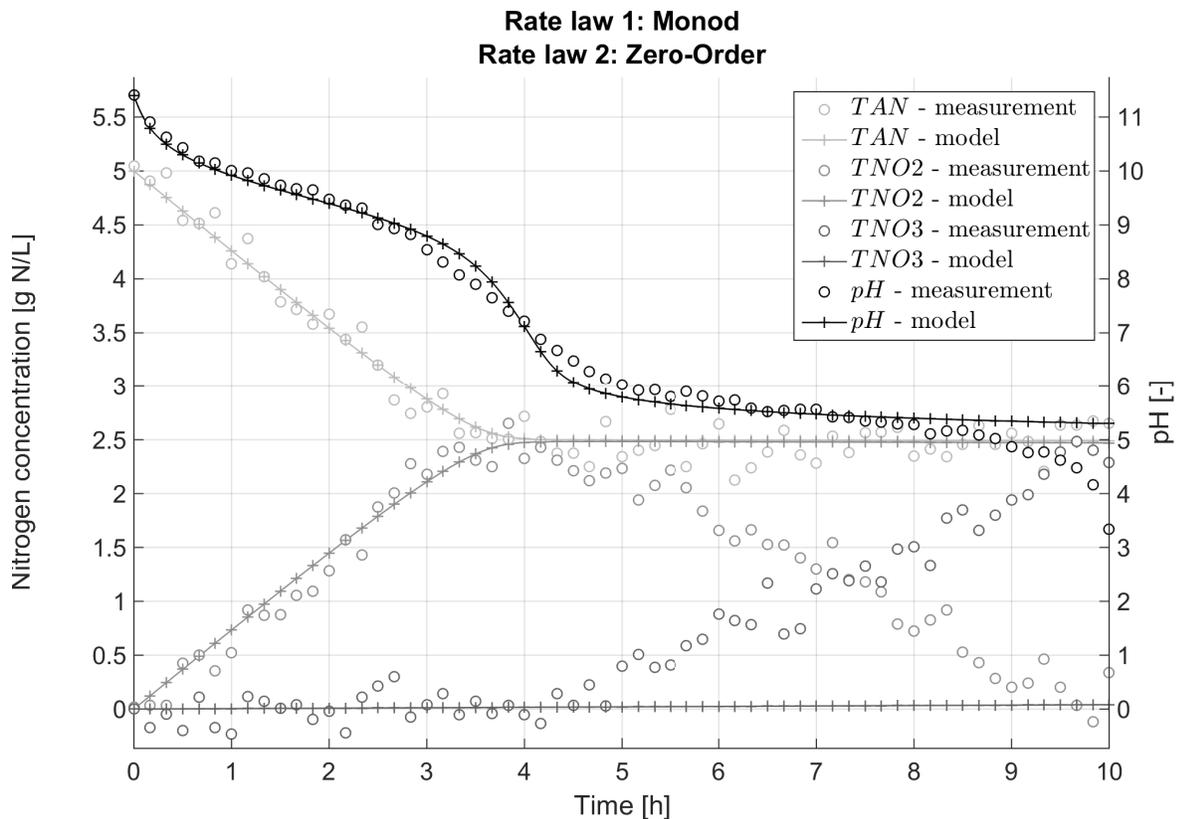


Figure S.5: **Method 1 - Simultaneous model identification - Model 3.** Measurements and simulation of the measured variables with model 3 after parameter estimation with the Nelder-Mead simplex method.

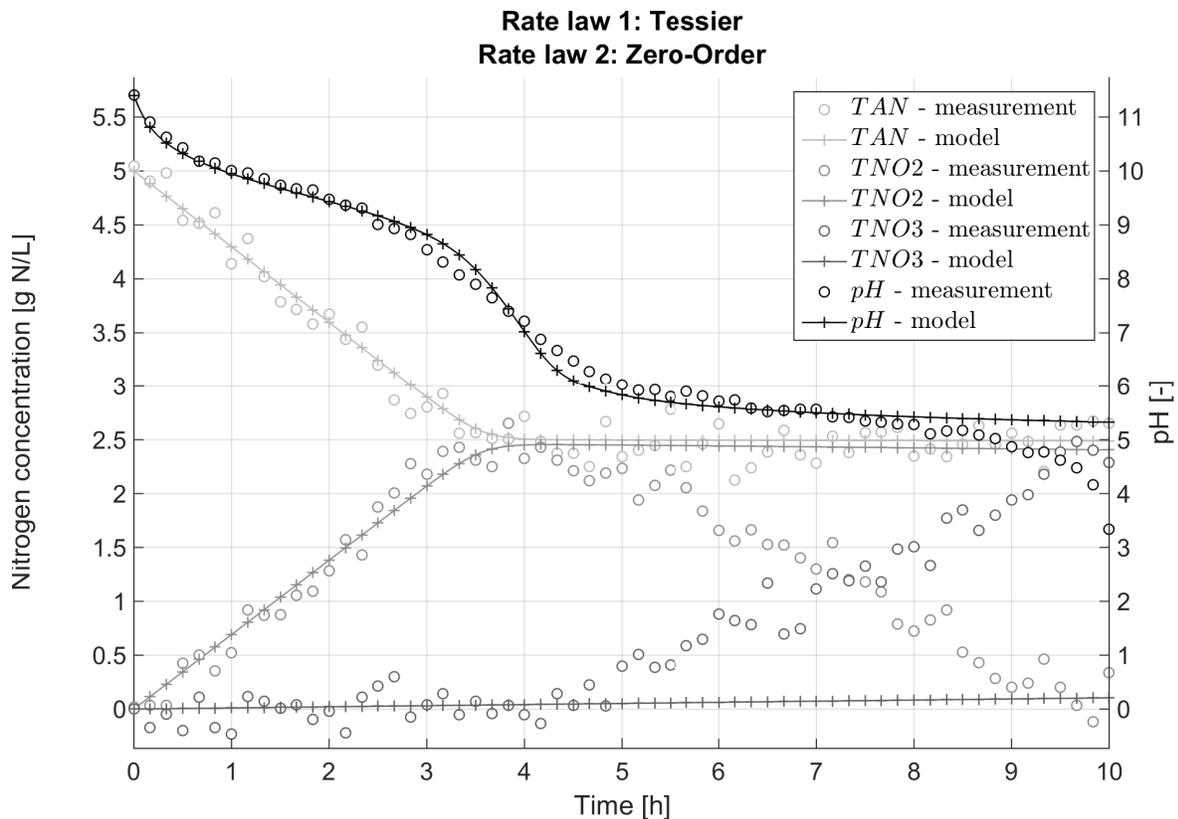


Figure S.6: **Method 1 - Simultaneous model identification - Model 4.** Measurements and simulation of the measured variables with model 4 after parameter estimation with the Nelder-Mead simplex method.

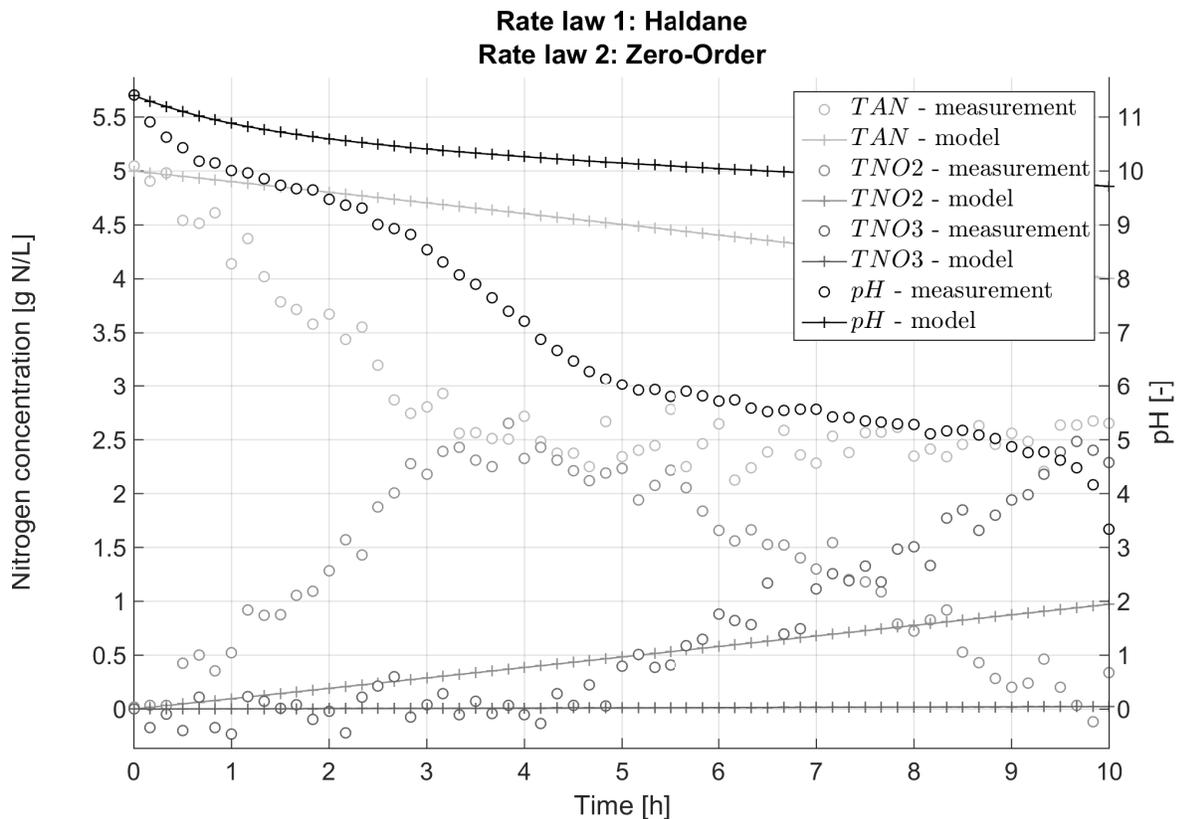


Figure S.7: **Method 1 - Simultaneous model identification - Model 5.** Measurements and simulation of the measured variables with model 5 after parameter estimation with the Nelder-Mead simplex method.

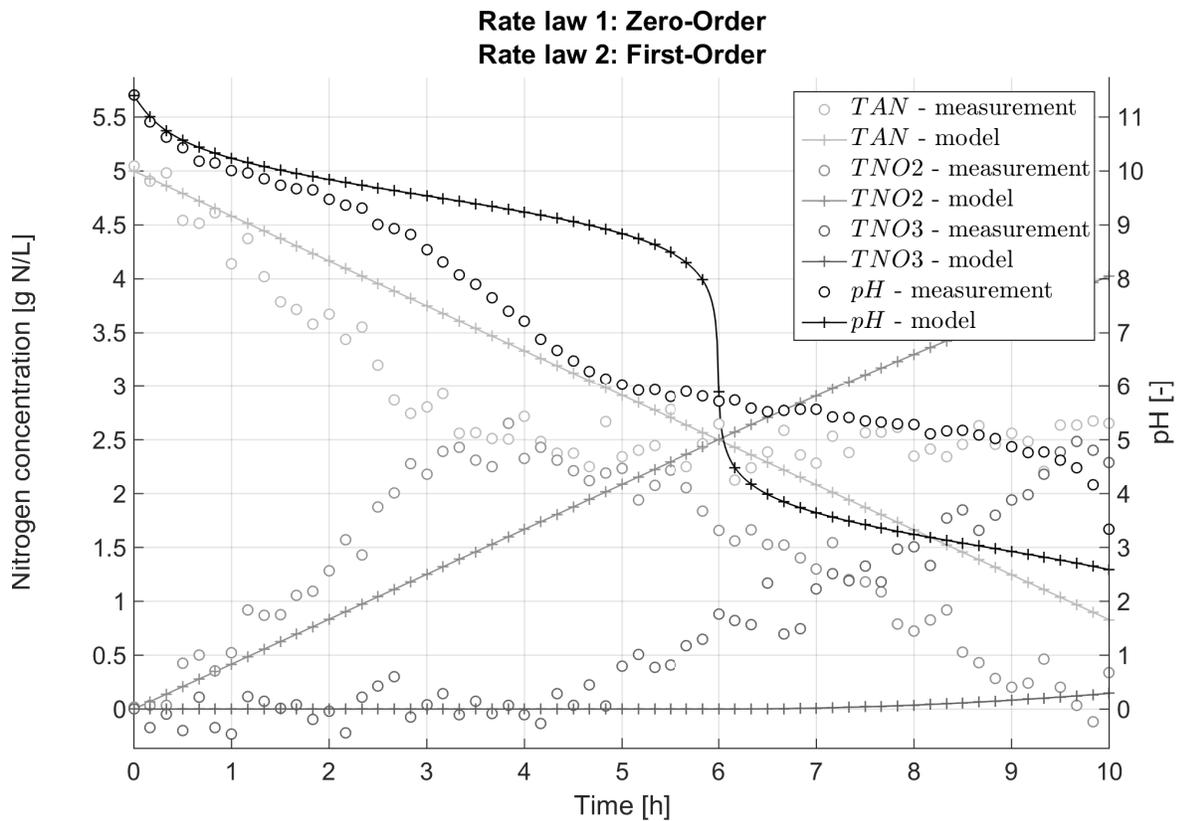


Figure S.8: **Method 1 - Simultaneous model identification - Model 6.** Measurements and simulation of the measured variables with model 6 after parameter estimation with the Nelder-Mead simplex method.

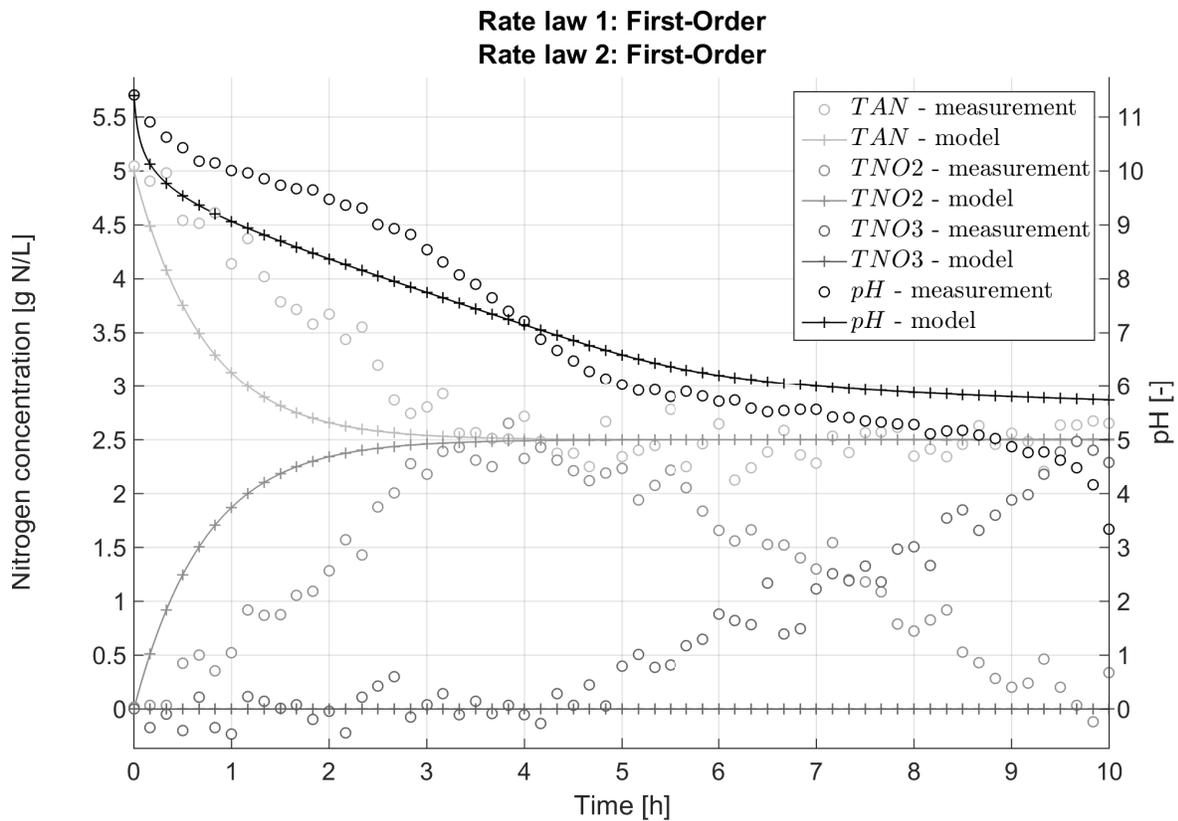


Figure S.9: **Method 1 - Simultaneous model identification - Model 7.** Measurements and simulation of the measured variables with model 7 after parameter estimation with the Nelder-Mead simplex method.

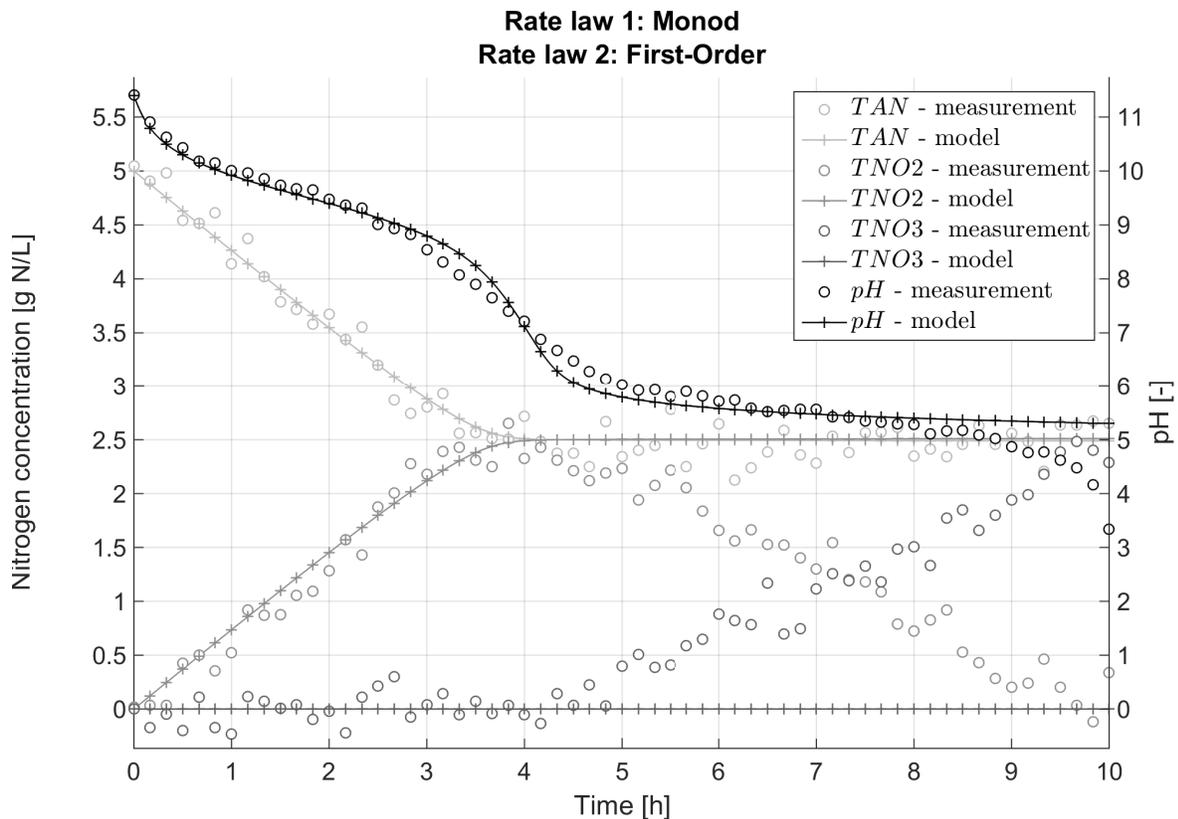


Figure S.10: **Method 1 - Simultaneous model identification - Model 8.** Measurements and simulation of the measured variables with model 8 after parameter estimation with the Nelder-Mead simplex method.

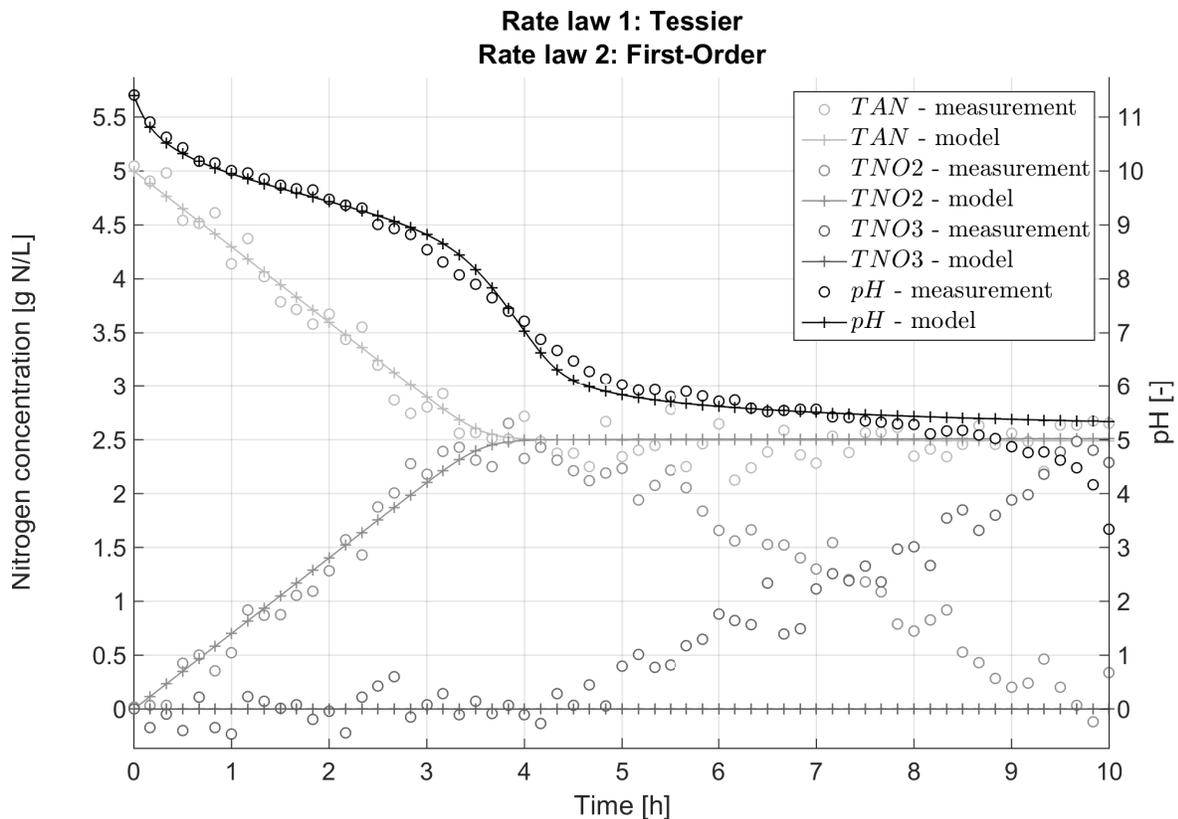


Figure S.11: **Method 1 - Simultaneous model identification - Model 9.** Measurements and simulation of the measured variables with model 9 after parameter estimation with the Nelder-Mead simplex method.

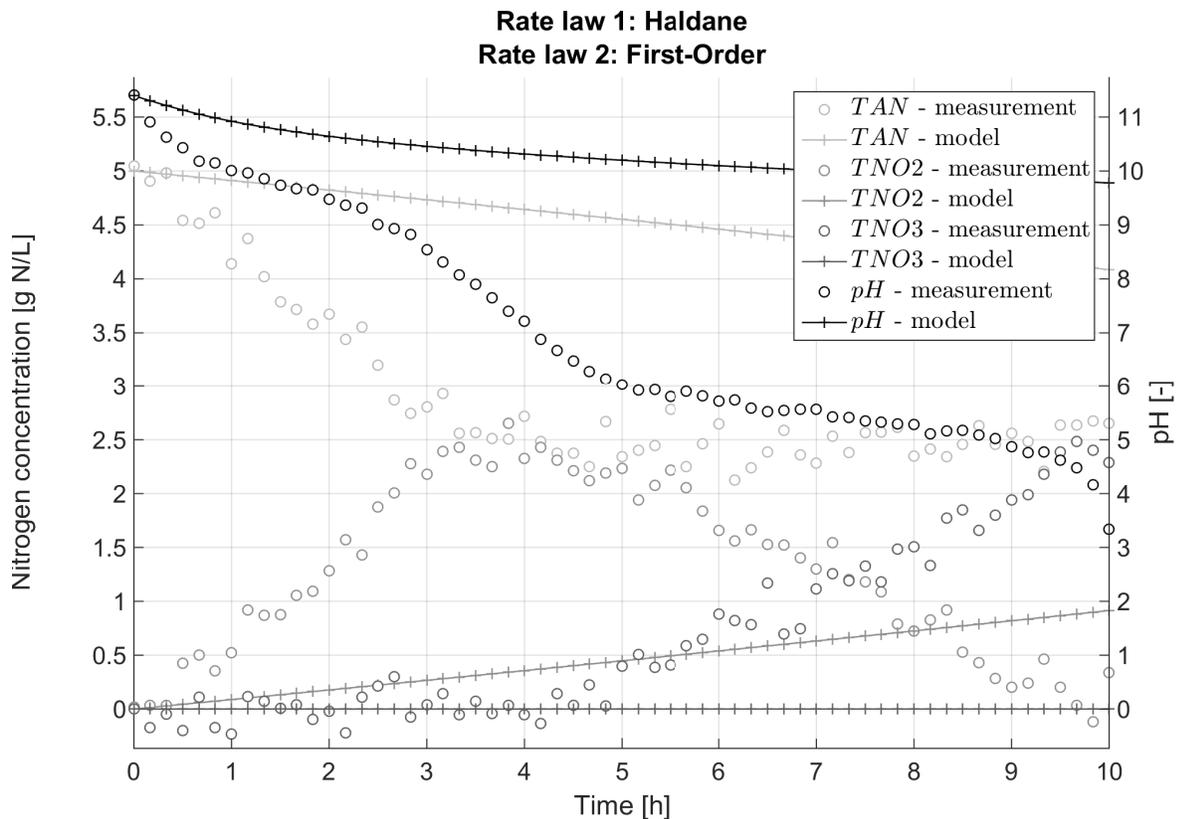


Figure S.12: **Method 1 - Simultaneous model identification - Model 10.** Measurements and simulation of the measured variables with model 10 after parameter estimation with the Nelder-Mead simplex method.

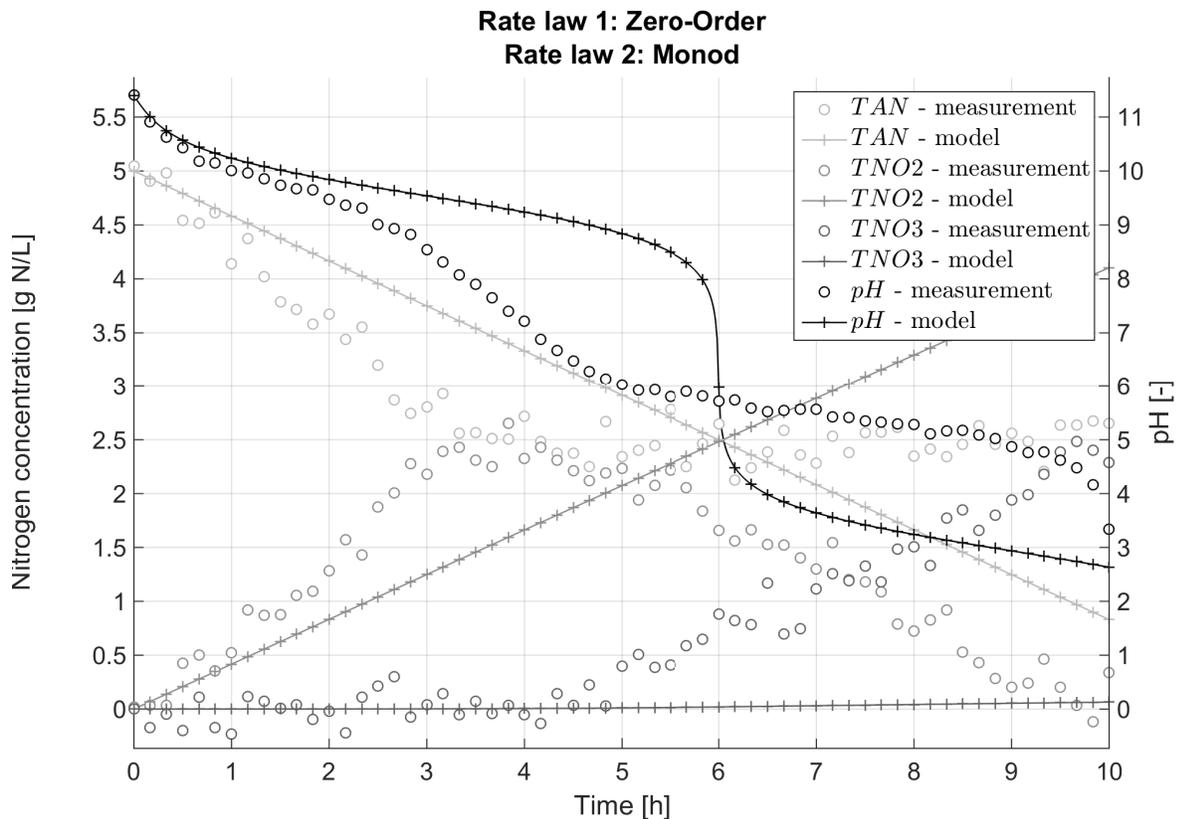


Figure S.13: **Method 1 - Simultaneous model identification - Model 11.** Measurements and simulation of the measured variables with model 11 after parameter estimation with the Nelder-Mead simplex method.

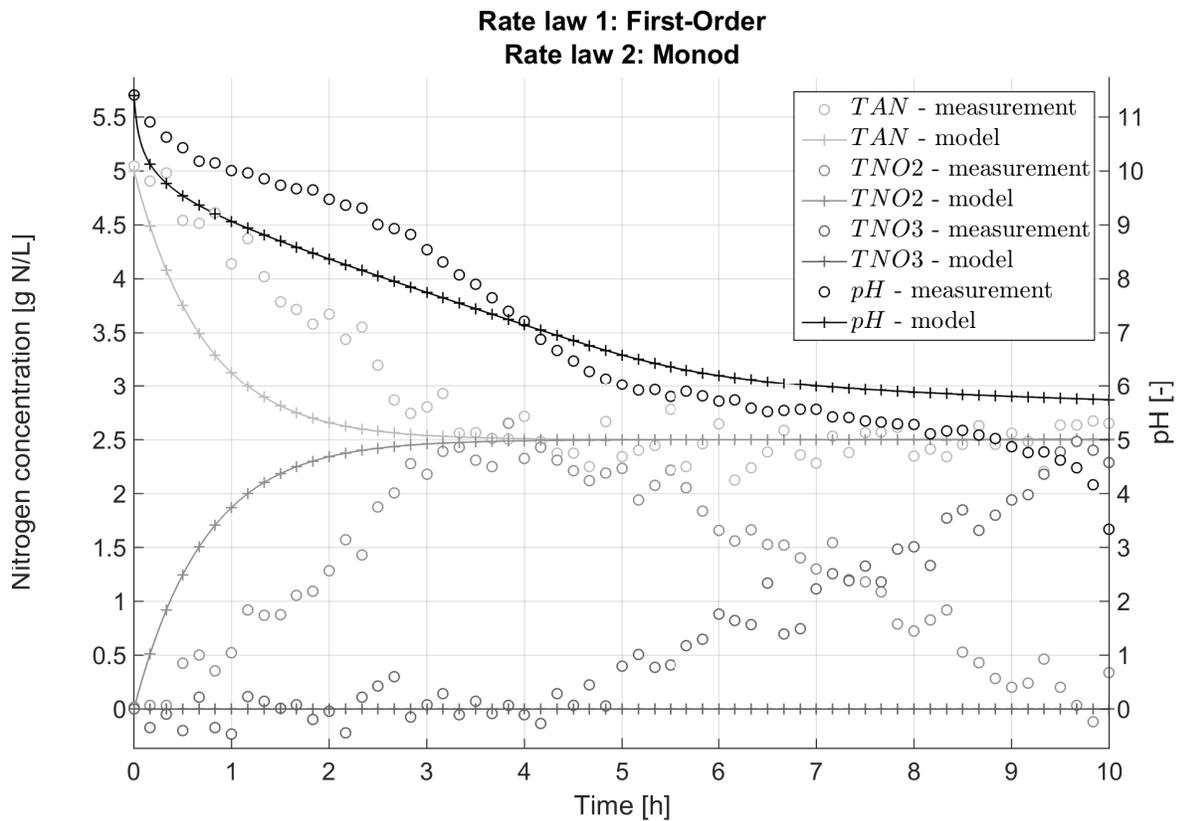


Figure S.14: **Method 1 - Simultaneous model identification - Model 12.** Measurements and simulation of the measured variables with model 12 after parameter estimation with the Nelder-Mead simplex method.

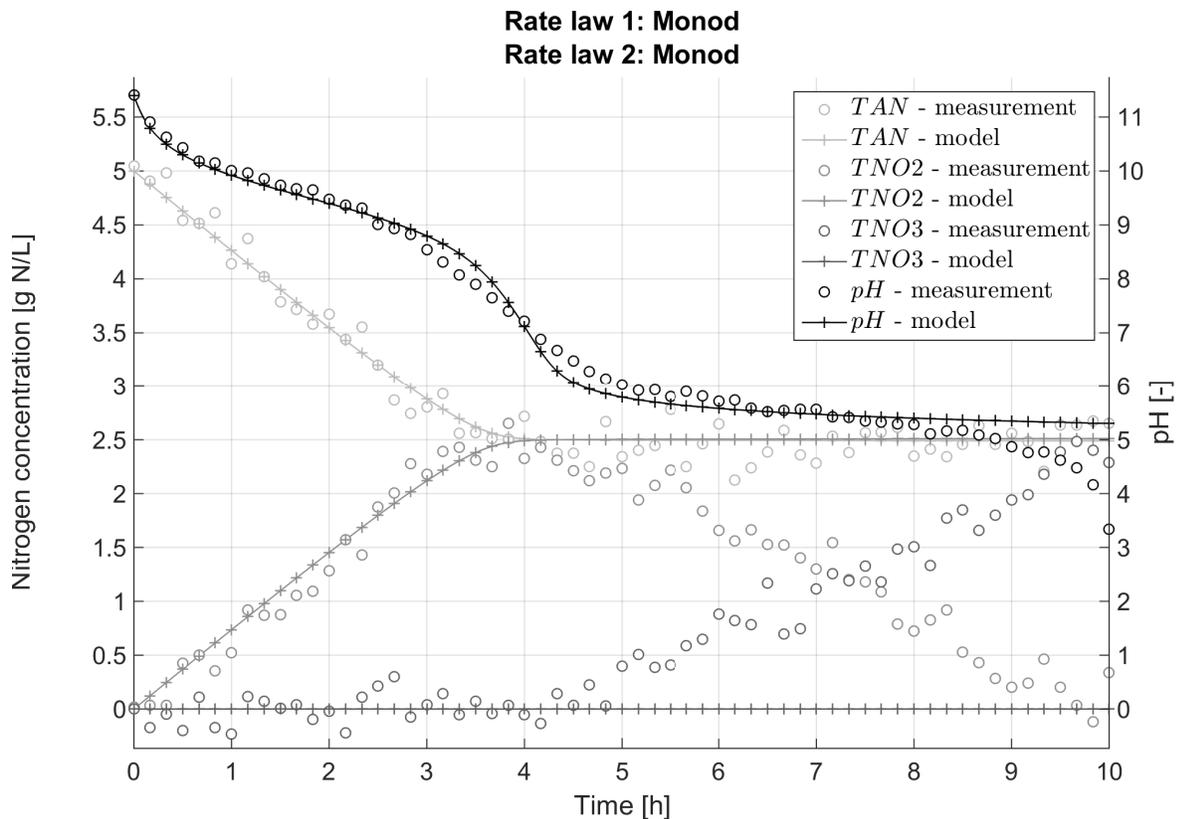


Figure S.15: **Method 1 - Simultaneous model identification - Model 13.** Measurements and simulation of the measured variables with model 13 after parameter estimation with the Nelder-Mead simplex method.

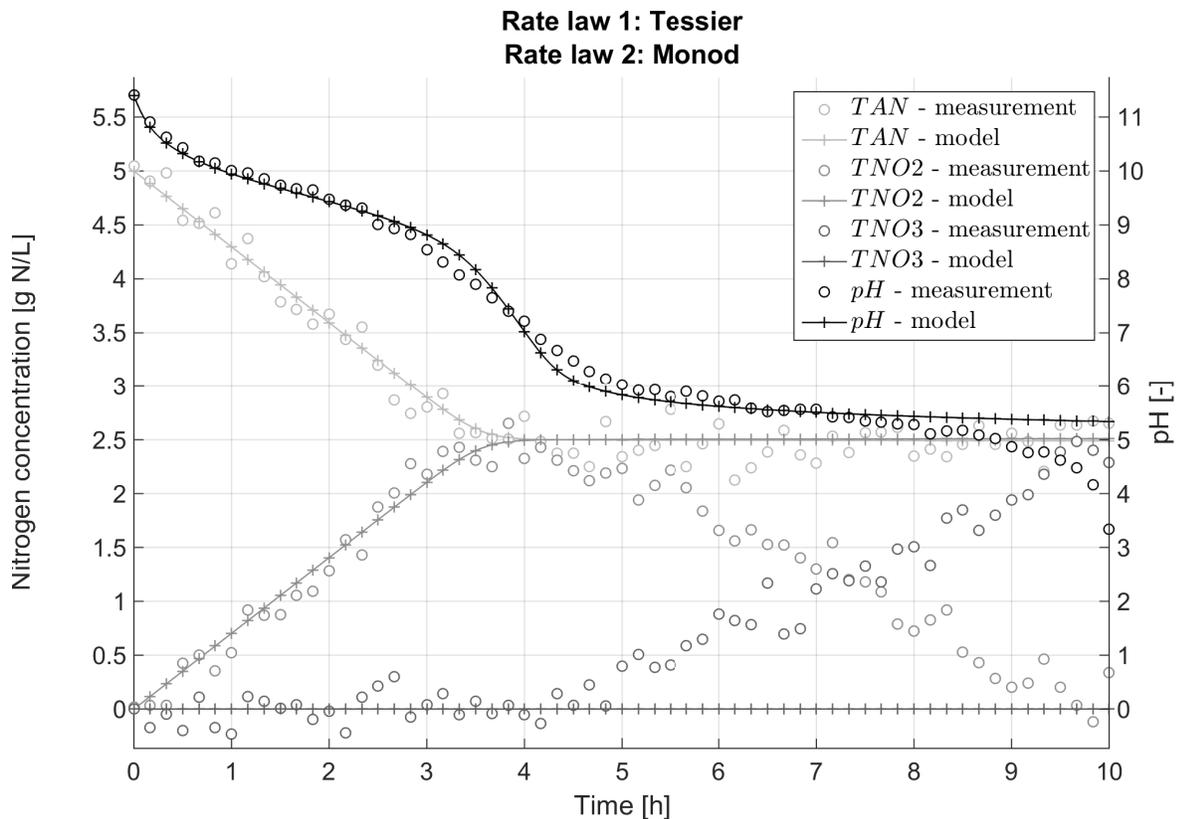


Figure S.16: **Method 1 - Simultaneous model identification - Model 14.** Measurements and simulation of the measured variables with model 14 after parameter estimation with the Nelder-Mead simplex method.

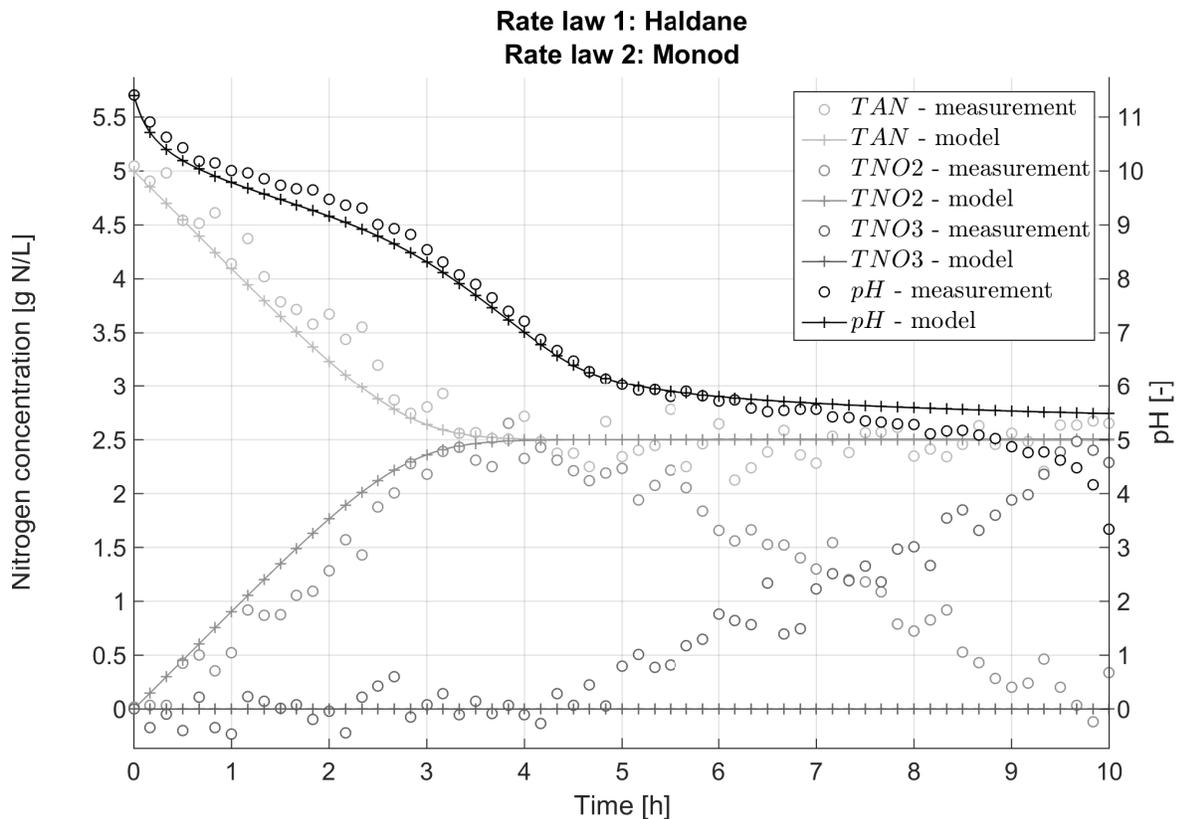


Figure S.17: **Method 1 - Simultaneous model identification - Model 15.** Measurements and simulation of the measured variables with model 15 after parameter estimation with the Nelder-Mead simplex method.

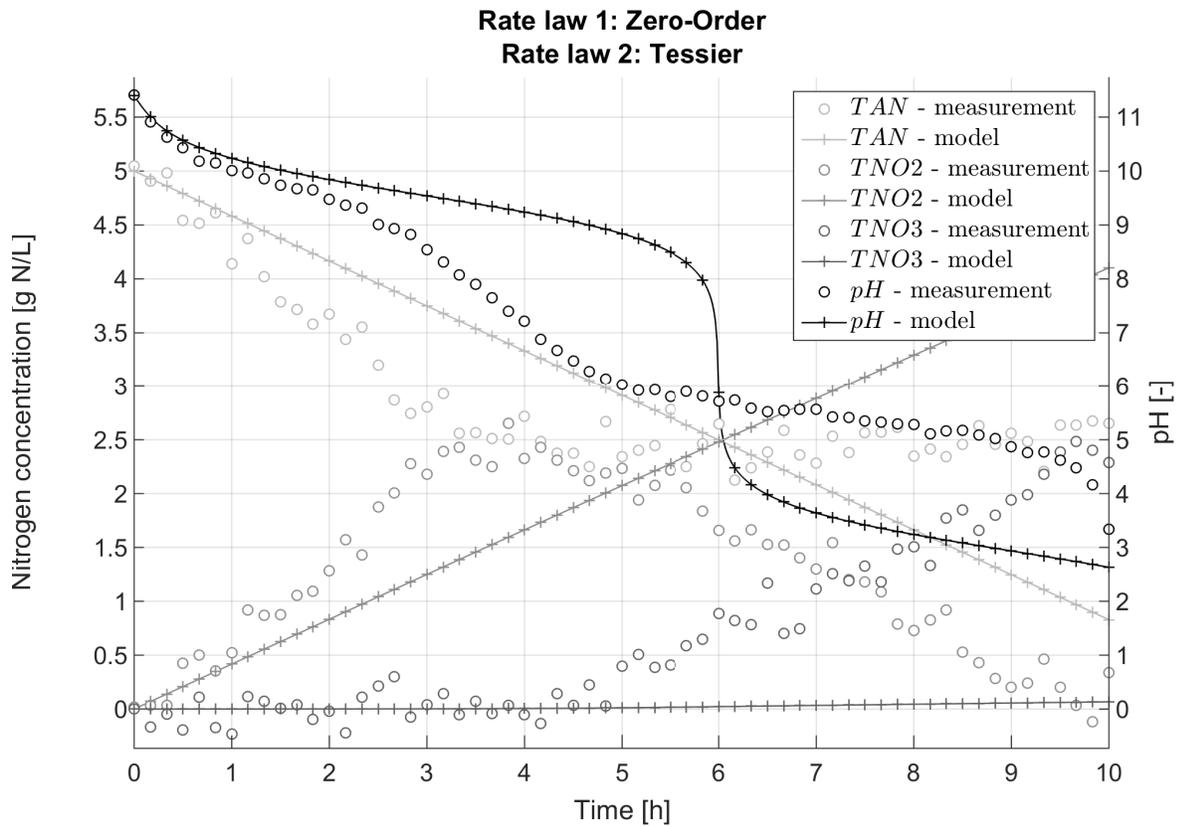


Figure S.18: **Method 1 - Simultaneous model identification - Model 16.** Measurements and simulation of the measured variables with model 16 after parameter estimation with the Nelder-Mead simplex method.

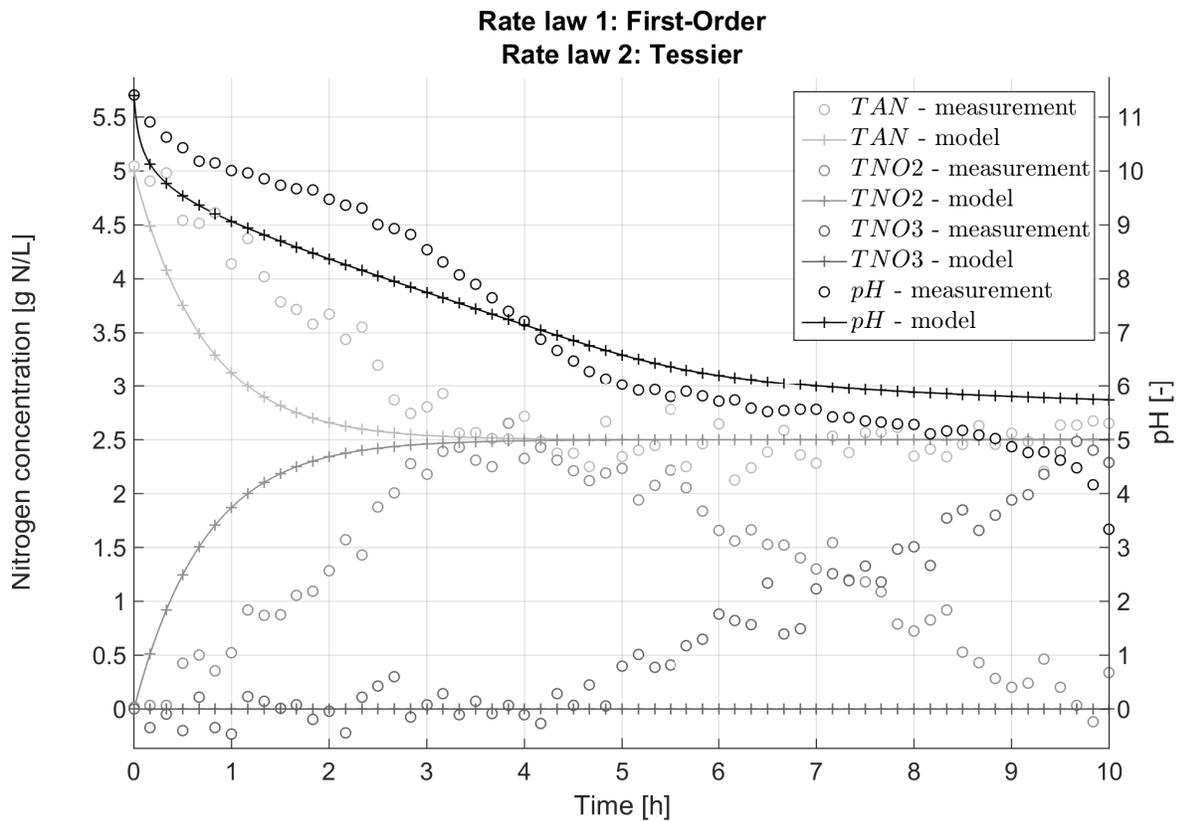


Figure S.19: **Method 1 - Simultaneous model identification - Model 17.** Measurements and simulation of the measured variables with model 17 after parameter estimation with the Nelder-Mead simplex method.

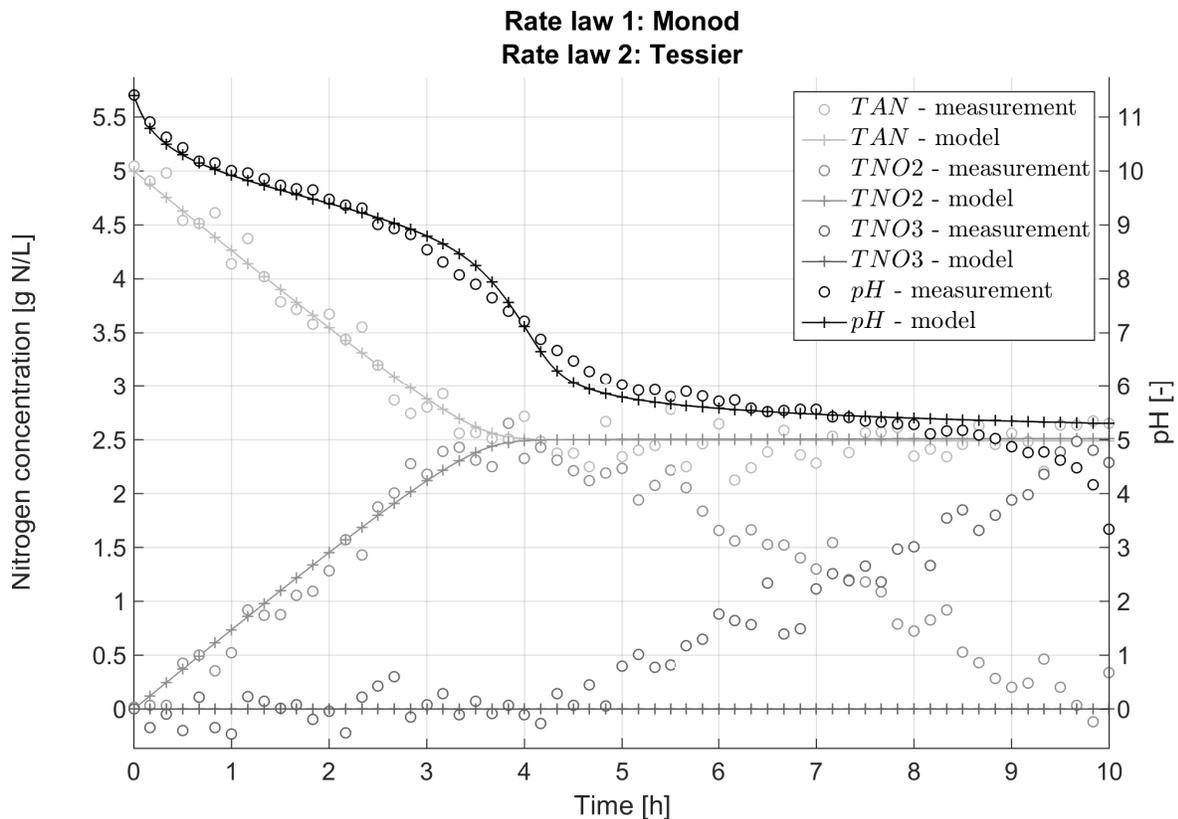


Figure S.20: **Method 1 - Simultaneous model identification - Model 18.** Measurements and simulation of the measured variables with model 18 after parameter estimation with the Nelder-Mead simplex method.

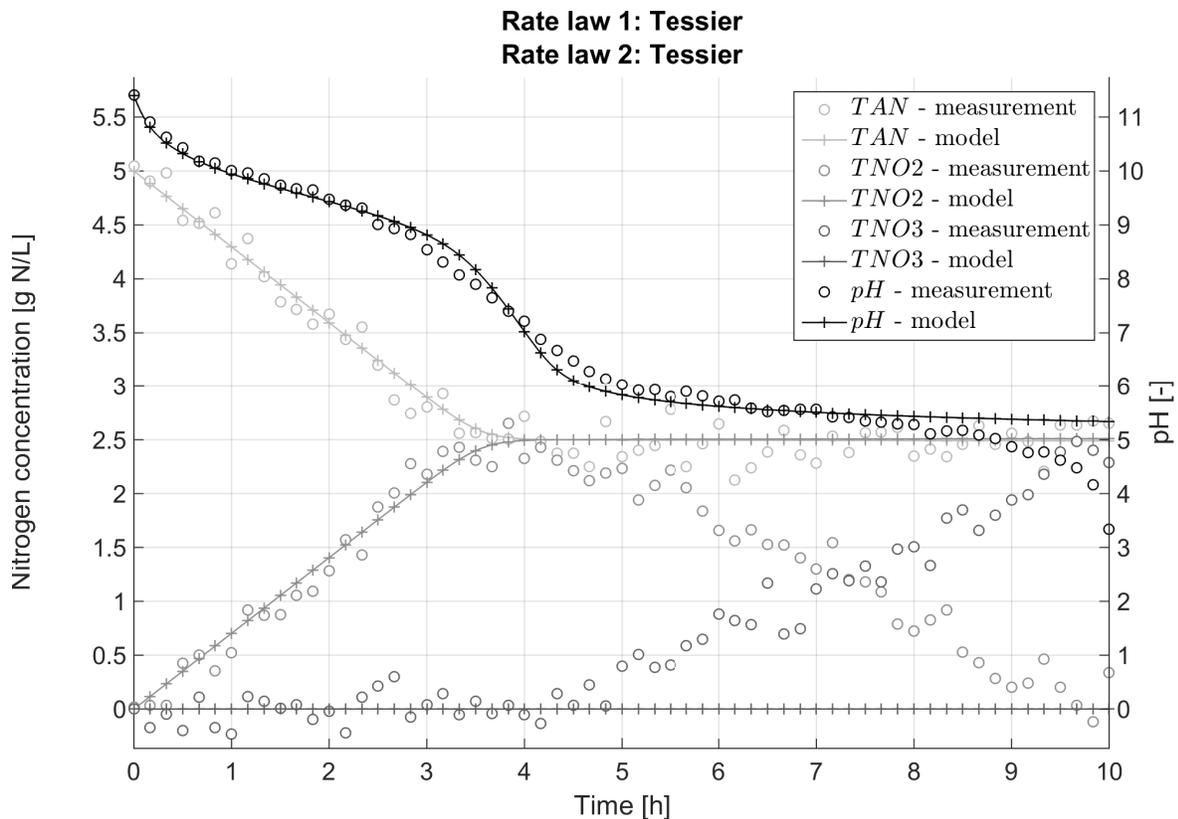


Figure S.21: **Method 1 - Simultaneous model identification - Model 19.** Measurements and simulation of the measured variables with model 19 after parameter estimation with the Nelder-Mead simplex method.

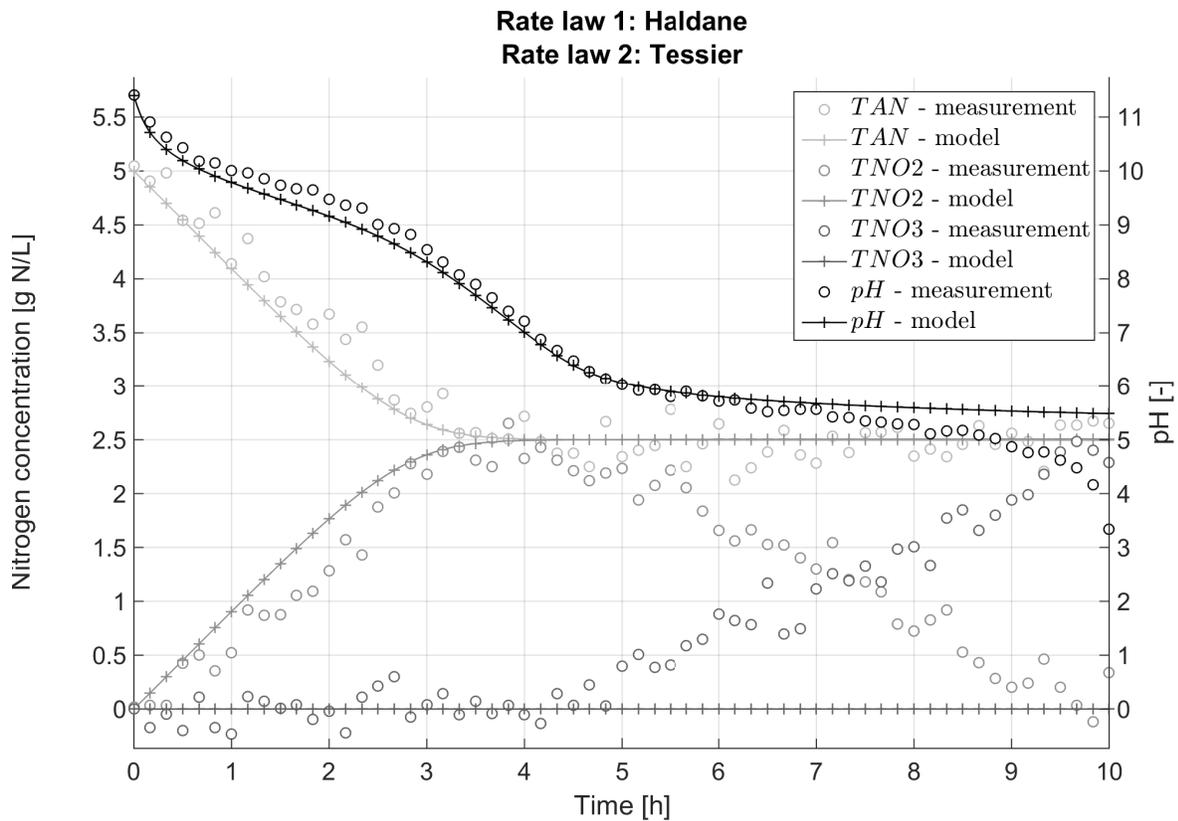


Figure S.22: **Method 1 - Simultaneous model identification - Model 20.** Measurements and simulation of the measured variables with model 20 after parameter estimation with the Nelder-Mead simplex method.

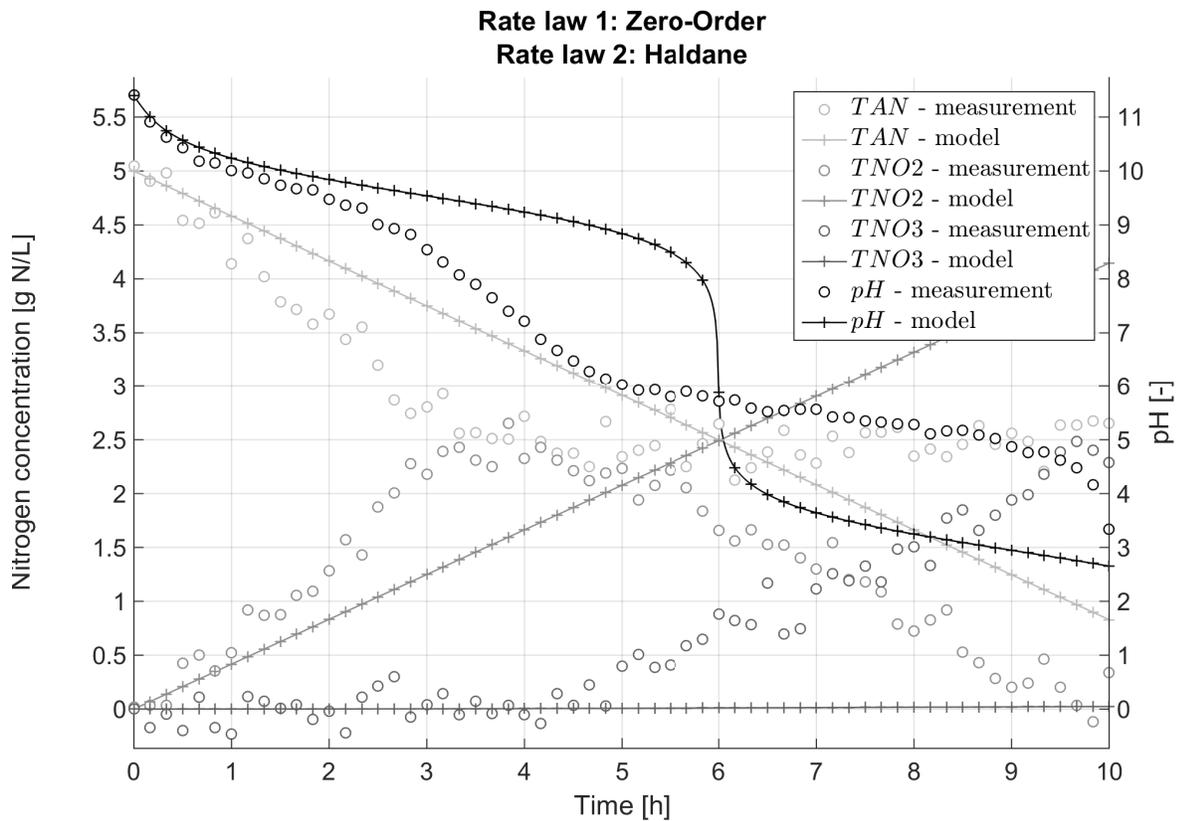


Figure S.23: **Method 1 - Simultaneous model identification - Model 21.** Measurements and simulation of the measured variables with model 21 after parameter estimation with the Nelder-Mead simplex method.

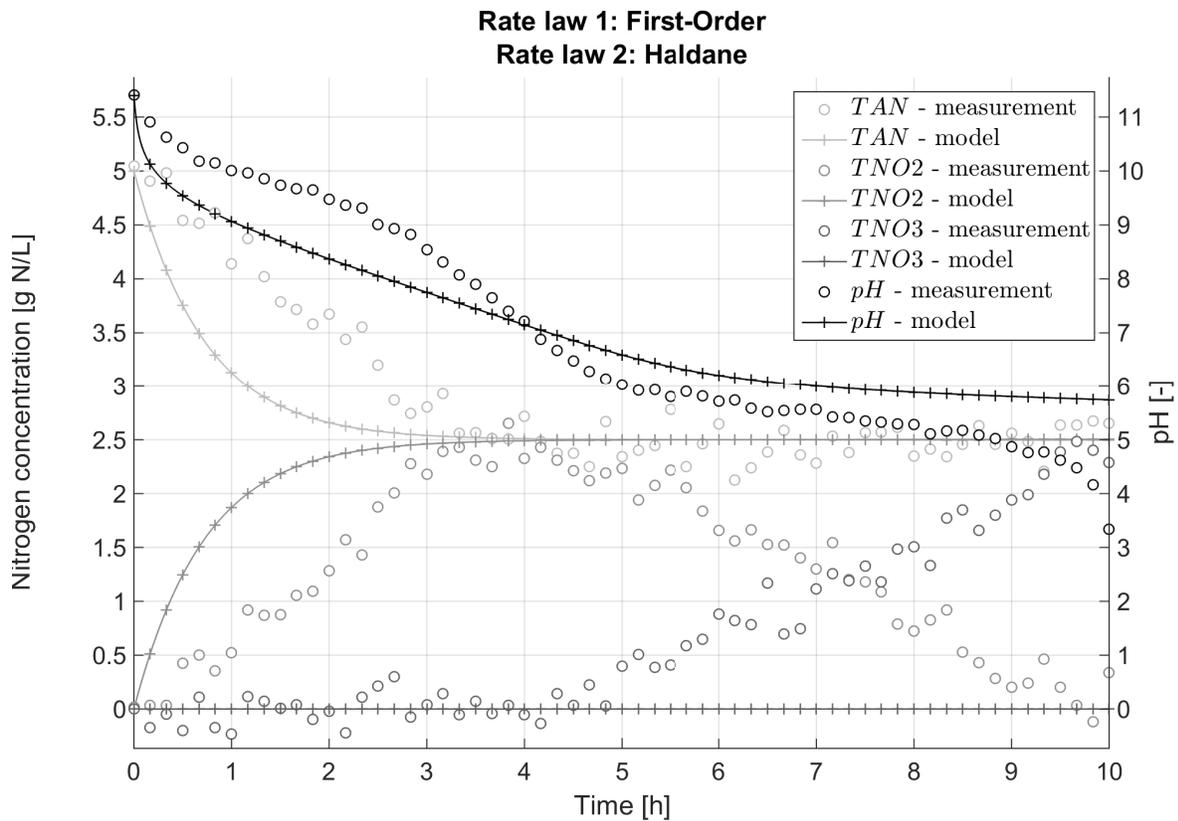


Figure S.24: **Method 1 - Simultaneous model identification - Model 22.** Measurements and simulation of the measured variables with model 22 after parameter estimation with the Nelder-Mead simplex method.

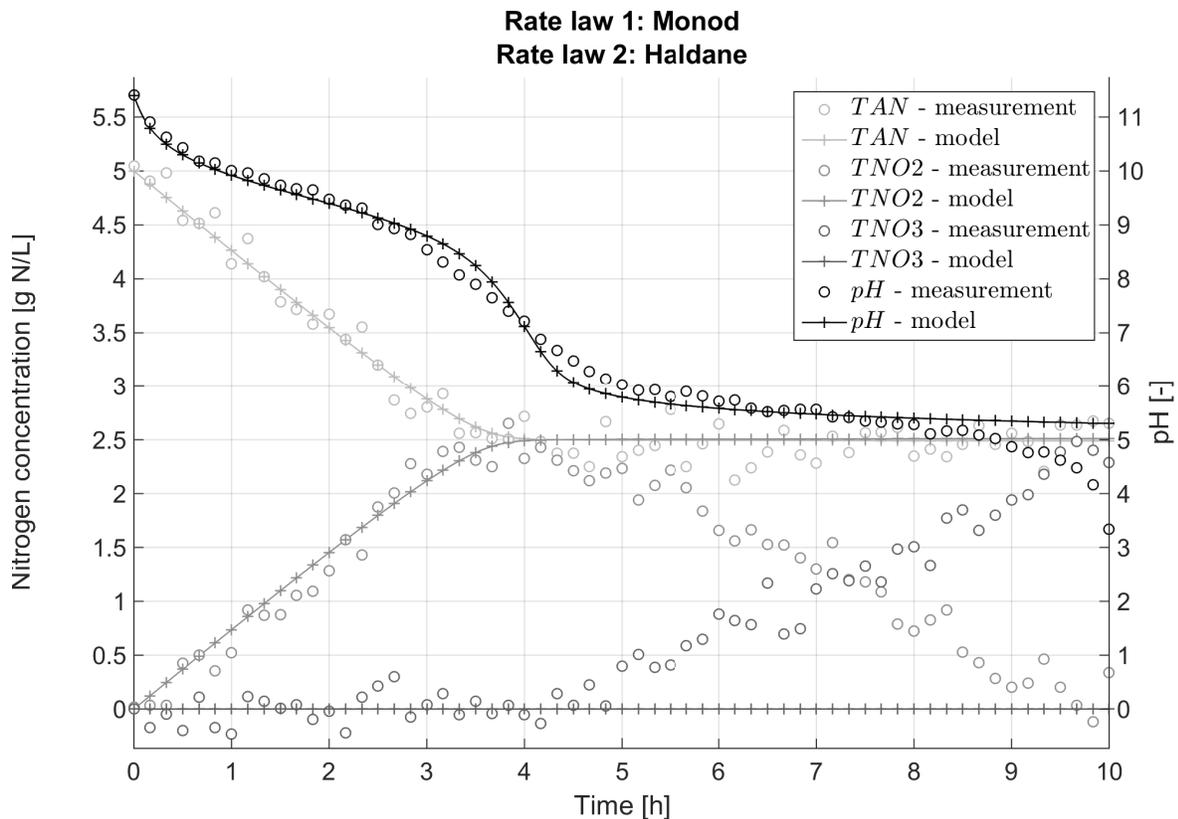


Figure S.25: **Method 1 - Simultaneous model identification - Model 23.** Measurements and simulation of the measured variables with model 23 after parameter estimation with the Nelder-Mead simplex method.

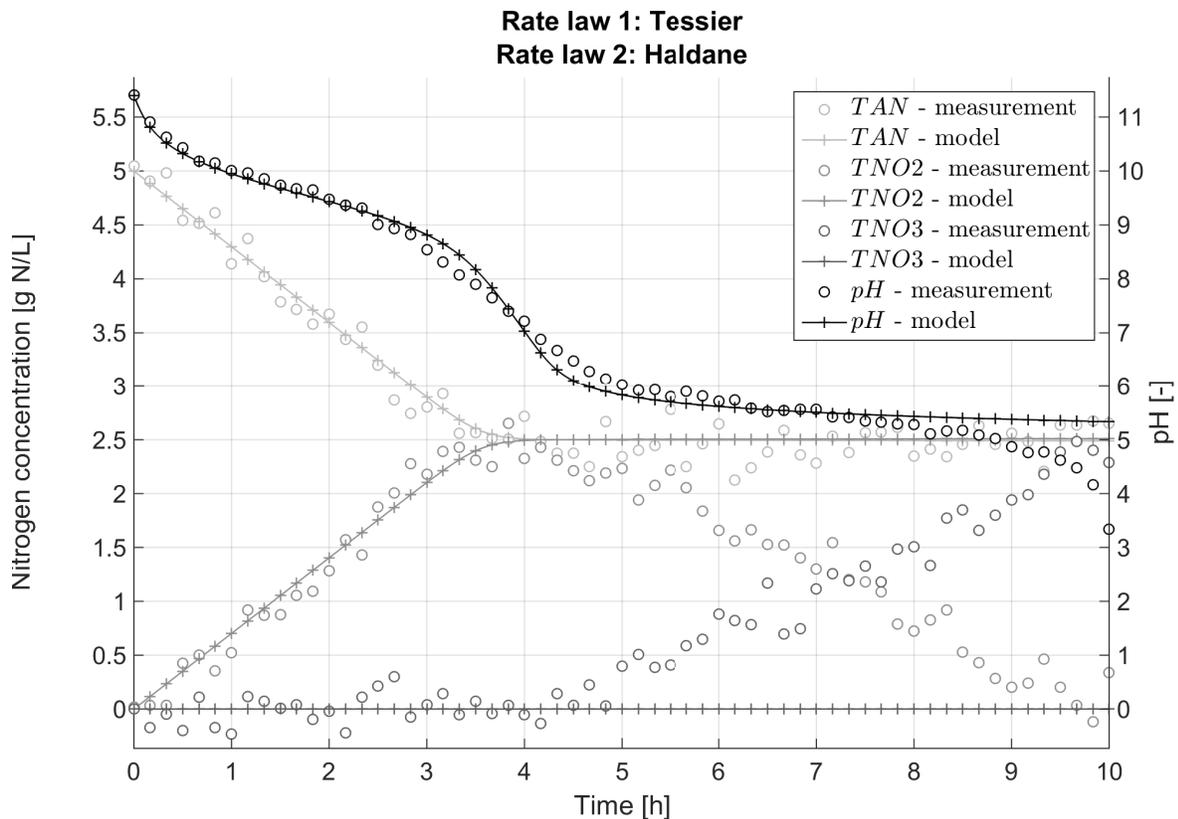


Figure S.26: **Method 1 - Simultaneous model identification - Model 24.** Measurements and simulation of the measured variables with model 24 after parameter estimation with the Nelder-Mead simplex method.

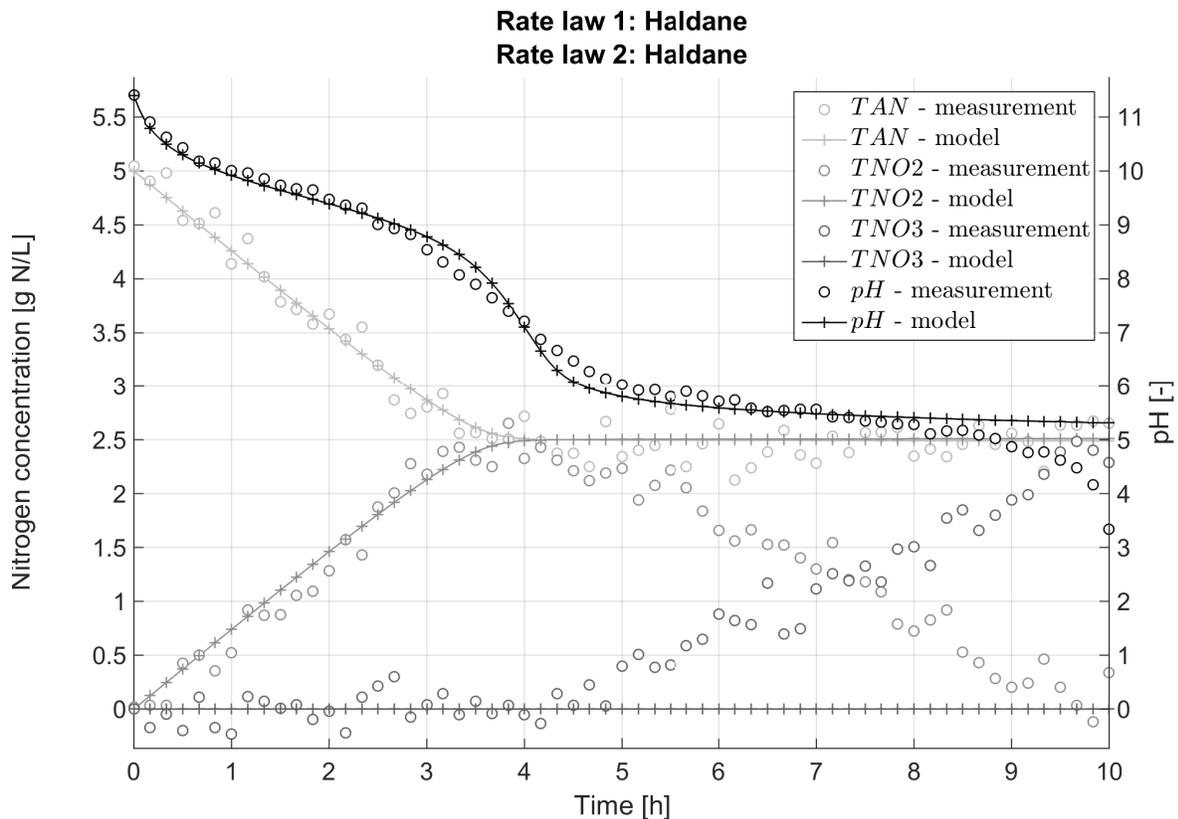


Figure S.27: **Method 1 - Simultaneous model identification - Model 25.** Measurements and simulation of the measured variables with model 25 after parameter estimation with the Nelder-Mead simplex method.

677 **Graphical TOC Entry**

678

