

1 Characterizing Long-term Wear and Tear of
2 Ion-Selective pH Sensors

3 Kito Ohmura^{a,b}, Christian M. Thürlimann^{a,c}, Marco Kipf^a, Juan Pablo
4 Carbajal^a, Kris Villez^{a,c,*}

5 ^a*Eawag, Überlandstrasse 133, CH-8600 Dübendorf, Switzerland*

6 ^b*Toshiba Infrastructure Systems & Solutions Corporation, Tokyo, Japan*

7 ^c*ETH Zürich, Institute of Environmental Engineering, 8093 Zürich, Switzerland*

8 **Abstract**

The development and validation of methods for fault detection and identification in wastewater treatment research today relies on two important assumptions: (i) that sensor faults appear at distinct times in different sensors and (ii) that any given sensor will function near-perfectly for a significant amount of time following installation. In this work, we show that such assumptions are unrealistic, at least for sensors built around an ion-selective measurement principle. Indeed, long-term exposure of sensors to treated wastewater shows that sensors exhibit important fault symptoms that appear simultaneously and with similar intensity. Consequently, our work suggests that focus of research on methods for fault detection and identification should be reoriented towards methods that do not rely on the assumptions mentioned above. This study also provides the very first empirically validated sensor fault model for wastewater treatment simulation and we recommend its use for effective benchmarking of both fault detection and identification

*Corresponding author, Email kris.villez@eawag.ch

methods and advanced control strategies. Finally, we evaluate the value of redundancy for the purpose of remote sensor validation in decentralized wastewater treatment systems.

9 *Keywords:* data quality, drift, fault detection and identification,
10 ion-selective electrodes, predictive maintenance, wastewater

11 **1. Introduction**

12 By several accounts, the lack of online sensor data quality poses a long-
13 standing challenge for both the advancement of environmental science and
14 engineering practice [17, 15, 18, 16, 9, 5]. It is therefore not surprising that
15 considerable time and energy has been invested in methods for automated
16 quality assessment and quality control of online measurement devices [e.g.,
17 23, 22, 6, 20, 1, 19, 30, 11].

18 Methods that are finding their way into practice today mainly consist of
19 sanity checks. In the authors' experience, these work rather well to detect
20 and classify a subset of commonly recognized fault symptoms, including out-
21 liers, spikes, stuck, and out-of-range values. For sensor faults that lead to
22 more subtle symptoms, current practice relies primarily on regular on-site
23 sensor maintenance, e.g. once every one or two weeks, to counter such subtle
24 faults. For unstaffed wastewater treatment plants, on-site maintenance may
25 be feasible economically only if this is limited to once per year. This practical
26 constraint to the adoption of quality assessment and control practices forms
27 the primary motivation for this study.

28 The literature suggests that data-analytical techniques can enable auto-
29 mated and remote detection of sensor faults. Without exception, such tech-
30 niques rely on redundant relationships and can therefore be categorized by
31 the type of redundancy that is used. A first category consists of techniques re-
32 lying on reference measurements and computing a deviation between online
33 sensor signal and the reference signal. A second category relies on hard-
34 ware redundancy by placing multiple online sensors, possibly built around
35 a distinct measurement principle, in the same location and then computing
36 deviations between them. A third category relies on temporal redundancy,
37 essentially assuming that meaningful changes in the sensor signal can only be
38 smooth when measured with a sufficiently high frequency. Finally, the fourth
39 category relies on spatial redundancy, relating signals produced at distinct
40 locations or for different measured variables. Examples of this last cate-
41 gory include both methods based on first principles, e.g. balance equations,
42 as well as methods rooted in statistical practice, e.g. principal component
43 analysis. Importantly, each of these advanced methods require tuning to
44 maximize the number of true alarms and to ensure suitable quality control
45 efforts while simultaneously minimizing the number of false alarms and fu-
46 tile maintenance actions. Invariably, such tuning is obtained by means of
47 a historical, fault-free data set from which acceptable limits for computed
48 residuals are derived. Consequently, this means that these methods rely on
49 the availability of representative data of an acceptable quality. In addition,
50 the use of most techniques implies that sensor fault symptoms can be as-

51 sumed to appear independently from each other, i.e. the probability that
52 two faults start at the same time is assumed to equal zero.

53 The prevalence of faults in actuators, sensors, and processes as well as
54 the complexity of the fault detection and identification (FDI) task, has led
55 to a plethora of methods that exploit one or more of the types of redundancy
56 discussed above. In fact, the wealth of literature as well as the number of
57 reviews on this or related topics [29, 27, 28, 9, 5] suggest that the science
58 and practice of FDI is all but settled, an observation also supported by no
59 free lunch theorems [33].

60 Despite the tremendous amount of research on FDI methods, little is ac-
61 tually known about the cause-and-effect relationships between sensor ageing,
62 the occurrence of sensor faults and failures, and the production of faulty data.
63 This is explained by the fact that the availability of information describing
64 the exact circumstances under which faults occur or faulty data is produced,
65 i.e. meta-data, is usually severely limited. This is the secondary motivation
66 of this study.

67 To facilitate performance evaluation of FDI tools, the formulation of sim-
68 ulation benchmarks has been an accepted practice in engineering sciences
69 [2, 7]. Similarly, the Benchmark Simulation Model No. 1 was conceived as
70 a way to test and compare innovative FDI and control strategies [10]. To-
71 day, it is primarily used as a starting point for a family of plant-wide models
72 of water resource recovery facilities [12, 31]. Actual benchmarking of FDI
73 methods has been limited to one study so far [6]. The BSM family includes

74 a set of sensor models which include sensor faults and this allows the user to
75 add realism to the sensor signals. The simulated sensor faults always start at
76 a time that is substantially later than the start of the simulated time. This
77 provides ideal conditions for FDI method tuning as high-quality sensor data
78 are always present in the first sections of the simulated data set. Moreover,
79 a simulated fault always appears independently of any other sensor fault, i.e.
80 no two sensor faults are simulated to start at the same time or with the same
81 direction or magnitude. We expect that the situation in real-world condi-
82 tions is very different. We thus hypothesize that typical fault symptoms will
83 appear at the same time and with similar directions and magnitudes when
84 exposed to the same harsh medium, especially when the same measurement
85 principle is applied. Evaluating the merit of this hypothesis is the tertiary
86 motivation of this study.

87 The following paragraphs are focused on the results and conclusions
88 drawn directly from experimental data obtained during a long-term sensor
89 exposure experiment. Additional insight is however obtained by studying a
90 variety of dynamic models to describe our measurements.

91 **2. Materials & Methods**

92 *2.1. Theoretical and real-world behavior of the ion-selective electrodes for pH* 93 *measurement*

94 The ion-selective measurement principle for pH measurement is under-
95 stood rather well. According to the Nernst equation [32] one measures an

96 electric potential E (in mV), which is related to the activity of the protons,
97 $[H^+]$, in the measured medium in steady state:

$$E = E^0 + \frac{RT}{F} \ln ([H^+]) \quad (1)$$

98 where E^0 is the reference potential, F is the Faraday constant [96485.33289 $C mol^{-1}$
99 21], $[H^+]$ is the proton activity in the reference cell, R is the molar gas con-
100 stant [8.3144598 $J mol^{-1} K^{-1}$ 21], and T is the temperature measured in
101 Kelvin. The pH is defined as $-\log [H^+]$ [3] so that $S(T)$ is the temperature-
102 specific sensitivity, which can be computed as:

$$S(T) = \frac{RT}{F \log(e)} \quad (2)$$

103 Most typically, pH sensors are designed to deliver 0 mV at pH 7 so that
104 E^0 is theoretically 0 mV. Similarly, the theoretical sensitivity at standard
105 ambient temperature and pressure (SATP) thus is $S(298.15) = 59.1593$ mV
106 per pH unit. Because the actual values of these parameters tend to deviate
107 from their theoretical values, it is common to identify their values through
108 a 2-point calibration procedure. At the engineering department at Eawag,
109 the most common practice is to use buffered calibration media with pH 4.01
110 and 7.00 for validation, followed by calibration when the absolute deviations
111 between the produced pH measurements and the known pH values exceed a

112 predetermined threshold. The data end user sets this threshold. Depending
113 on the application, this ranges from 0.1 to 0.4 pH units. The theoretical
114 potential at pH 4.01 and SATP is 177.0 mV.

115 2.2. Studied sensors

116 A total of 12 pH sensors are produced by Endress+Hauser (Reinach,
117 Switzerland). These sensors consist of 5 sensor types (T1-T5) whose exact
118 type cannot be revealed due to a confidentiality agreement. The first eight
119 sensors consist of pairs of four commercially available sensor types (T1-T4)
120 which are typically sold with a one-year warranty agreement. The first (sec-
121 ond) sensor in each pair is designated with an *a* (*b*), e.g. T1a, T1b. The last
122 4 pH sensors are replicates of a recently developed sensor prototype (T5) and
123 are referred to as T5a, T5b, T5c, and T5d.

124 The first three sensor pairs (T1-T3) have been in use throughout a long-
125 term exposure experiment which lasted for 731 days (Oct. 4th, 2016 – Oct.
126 4th, 2018). An overview of this experiment is given in Fig. 1. The 4th pair
127 (T4) has been in use during the first half year and was replaced with the
128 5th pair (T5) on April 3rd, 2017 (day 182) as (*i*) the T4 sensors exhibit a
129 long response time (not shown) and (*ii*) the opportunity arose to test the T5
130 prototypes. The T5a sensor stopped producing a meaningful signal on June
131 30th, 2017 (day 270) while T5b became faulty (details below) on August
132 31st, 2017 (day 332). These sensors were replaced with another sensor of
133 the same prototype (T5) on Oct. 2nd, 2017 (day 364). In this last pair, one

134 sensor (T5d) failed within 1 day (day 365) while the other (T5c) has been
135 fully functional until the end of the experiment.

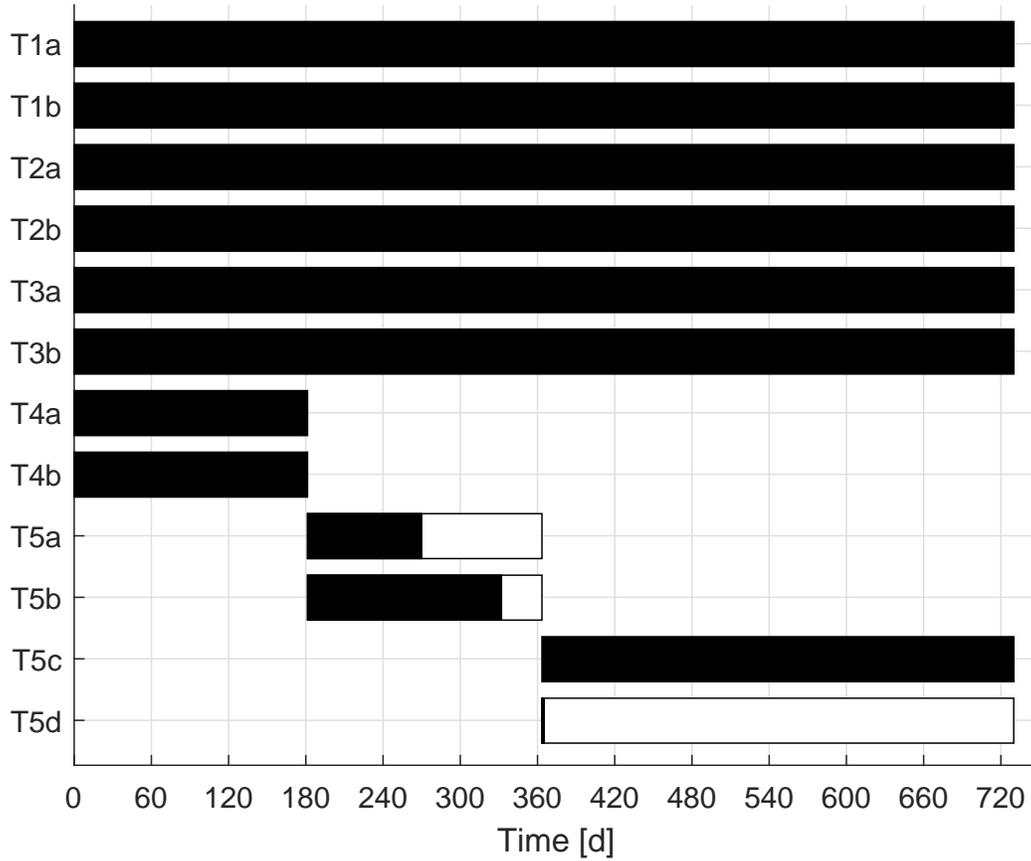


Figure 1: **Overview of the complete experimental campaign.** The periods of sensor exposure are indicated by rectangles. The periods during which the sensors produced meaningful data are marked black.

136 *2.3. Long-term exposure experiment*

137 The sensors are exposed to the contents of a reactor used primarily to
138 study advanced control strategies for nitrite accumulation prevention in a

139 urine nitrification process [26]. To this end, the nitrified urine is pumped
140 through a closed tube made from PVC with a flow rate of 43 L/h. The design
141 of this tube equipped with sensor-holding locks is shown in the *Supplementary*
142 *Information (Section B)*.

143 The treated urine is from anthropogenic origin during the whole experi-
144 mental period. The treated urine was collected from male lavatories in the
145 Forum Chriesbach building at Eawag, with exception of the period from day
146 April 30th, 2018 to June 21st, 2018 (day 574-625), when it was collected from
147 female lavatories in the same building. From October 4th, 2017 to November
148 24th, 2017 (day 366 to 417), the reactor was additionally fed with a nitrite
149 stock solution. During the experimental period, the measured concentra-
150 tions of nitrogen species in the nitrified urine ranged between 1180 and 2730
151 mgN/L (mg atomic nitrogen per liter) for total ammonia, 0 and 82 mgN/L
152 for nitrite, and 1290 and 2720 mg/L for nitrate. These measurements are
153 copied from [26] and are shown in the *Supplementary Information (Section*
154 *C)*. The pH value of the nitrified urine, as measured by two independent
155 and regularly calibrated pH sensors installed directly in the reactor, ranged
156 between 5.7 and 7.3.

157 2.4. Sensor characterization tests

158 At regular intervals, the sensors were removed from their normal position
159 and exposed to other media for sensor characterization. This was executed
160 47 times in total. The exact times of these sensor characterization tests are

161 listed in the *Supplementary Information (Section G.1)*. Two pairs of tests
 162 were executed on the same day to ensure acceptable experimental repro-
 163 ducibility (day 70: tests 11-12; day 351: tests 29-30). The selected media
 164 include (C4) pH 4.01 calibration solution (CPY20-C10A1, Endress+Hauser,
 165 Reinach, Switzerland); (C7) pH 7.00 calibration solution (CPY20-E10A1,
 166 Endress+Hauser, Reinach, Switzerland); (U4) nitrified urine at pH 4; (U7)
 167 nitrified urine at pH 7; and (W) tap water. For the present work, only the
 168 exposure to W, C4, and C7 is relevant. This occurs in five distinct phases
 169 (P0-P4), each lasting at least 5 minutes and exposing the sensors to W, C4,
 170 C7, C4, and W in this order. Exemplary results are shown in Fig. 2 and
 171 discussed in detail below.

172 Raw potential measurements recorded during P1, P2, and P3 are used
 173 to compute the offset (\tilde{E}^0) and two measurements of the sensitivity (\tilde{S}_D and
 174 \tilde{S}_R). To this end, the following steps are applied for every sensor and every
 175 sensor characterization test [4]:

- 176 1. Compute the median value among the potential measurements collected
 177 in P1, P2, and P3 between 2 and 1 minutes before the start of the next
 178 phase (P2, P3, and P4). Refer to these values as E^{P1} , E^{P2} , and E^{P3}
- 179 2. The sensor offset is defined as $\tilde{E}^0 = \tilde{E}^{P2}$.
- 180 3. The decay potential sensitivity is defined as $\tilde{S}_D = \frac{\tilde{E}^{P1} - \tilde{E}^{P2}}{7.00 - 4.01} = \frac{\tilde{E}^{P1} - \tilde{E}^{P2}}{2.99}$.
- 181 4. The decay potential sensitivity is defined as $\tilde{S}_R = \frac{\tilde{E}^{P3} - \tilde{E}^{P2}}{7.00 - 4.01} = \frac{\tilde{E}^{P3} - \tilde{E}^{P2}}{2.99}$.

182 These steps are demonstrated below with a practical example.

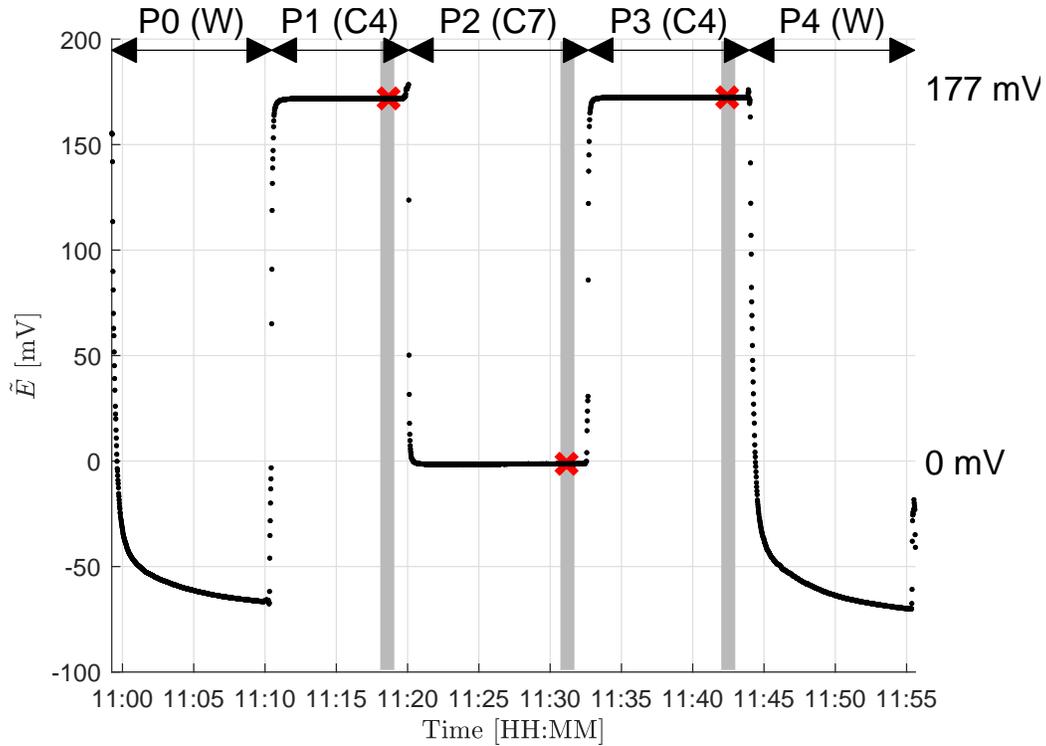


Figure 2: **Exemplary sensor characterization test.** Raw data obtained in the first sensor characterization test with sensor T1a. The measured potential decays during P0, P2, and P4, while it increases during P1 and P3. Steady state is reached quickly in P1, P2, and P3. The theoretical potential values for P1, P2, and P3 are indicated with dashed horizontal lines. Grey shading indicates the data used to obtain the potential measurements (2 to 1 minute before phase change). The selected median potential values are shown with red crosses.

183 *2.5. Drift model*

184 The results shown below indicate that the offset significantly varies over
 185 time while the sensitivity remains remarkably stable in all studied sensors.

186 We describe the observed drift of the offset by means of two models.

187 *2.5.1. Model 1 - Constant trend followed by linear trend*

188 For the first model, we apply a modified version of the excessive drift
 189 model proposed for the BSM family [18]. This model simulates $E^0(t)$, the
 190 sensor offset, as:

$$E^0(t) = d_o + r_d H(t - t_f) \quad (3)$$

191 with d_o the initial offset, r_d the drift rate parameter, $H(\cdot)$ the Heaviside
 192 function ($H(a) = 1$ if $a \geq 0$, $H(a) = 0$ otherwise), t the time since sensor
 193 installation, and t_f the time of the drift onset. The applied modification con-
 194 sists of adding the parameter d_o . To fit this model, the offset measurements,
 195 $\tilde{E}^0(t_h)$, collected at discrete time instants t_h , are assumed to exhibit inde-
 196 pendently and identically distributed measurement errors, ϵ_h , drawn from a
 197 normal distribution with zero mean and standard deviation, σ_ϵ :

$$\tilde{E}^0(t_h) = E^0(t_h) + \epsilon_h, \epsilon_h \sim N(0, \sigma_\epsilon) \quad (4)$$

198 Values for the 4 parameters d_o , t_f , r_d , and σ_ϵ are obtained independently
 199 for all sensors through maximum likelihood estimation (MLE). Once cali-
 200 brated, the model is used to obtain the estimated mean and point-wise stan-
 201 dard deviations for the sensor offset, $\mu_1(t) = \mathbb{E}(E^0(t))$ and $\sigma_1(t)$, while using
 202 the estimates of t_f and σ_ϵ as fixed hyperparameter values.

203 *2.5.2. Model 2 - Integrated Brownian motion for a single sensor*

204 In model 2, we assume instead that the recorded offset measurements
 205 are generated by an integrated Brownian motion. This is a continuous-time
 206 stochastic process, which reflects that the drift rate is subject to unmeasured
 207 disturbances:

$$\dot{r}_d(t) = \gamma(t)dt, r_d(0) = r_{d,o}, \gamma(t) \sim N(0, \sigma_\gamma), \quad (5)$$

$$\dot{E}^0(t) = r_d(t)dt, E(0) = d_o, \quad (6)$$

$$\tilde{E}^0(t_h) = E^0(t_h) + \epsilon_h, \epsilon_h \sim N(0, \sigma_\epsilon) \quad (7)$$

208 This model also includes 4 parameters: the initial drift rate ($r_{d,o}$); the
 209 initial offset (d_o); an input noise standard deviation controlling the rate by
 210 which the drift rate changes (σ); and an output noise standard deviation
 211 (σ_ϵ). As with model 1, parameter values are obtained through MLE. This is
 212 achieved by formulating the above process as a Gaussian process [14]. This
 213 also enables to compute expected values and associated point-wise standard
 214 deviations, $\mu_2(t) = \mathbb{E}(E^0(t))$ and $\sigma_2(t)$, with the estimates of σ_γ and σ_ϵ now
 215 used as fixed hyperparameter values.

216 *2.5.3. Model 3 - Integrated Brownian motion for multiple sensors*

217 A third model is derived from Eqs. 5-7 by considering that two sensors
 218 of the same type may be characterized by distinct initial conditions ($r_{d,o}$,

219 d_o) but the same noise parameters $(\sigma_\epsilon, \sigma_\gamma)$. This lead to a model with six
220 parameters $(d_o^a, d_o^b, r_{d,o}^a, r_{d,o}^b, \sigma_\epsilon, \sigma_\gamma)$, instead of two models with 4 parameters
221 each. Their values are again obtained via MLE and used to obtain calibrated
222 predictions $(\mu_3(t) = \mathbb{E}(E^0(t)), \sigma_3(t))$, once again using the estimates of σ_γ
223 and σ_ϵ as fixed hyperparameter values.

224 2.5.4. Model evaluation

225 The proposed models are evaluated through visual inspection of the mea-
226 surements, predictions, and residuals between the measurements and predic-
227 tions. In the present case, such a visual inspection is considered sufficient to
228 select a suitable model.

229 2.5.5. Implementation

230 All data collected during the sensor characterization tests and all code
231 necessary to reproduce our results is added in the *Supplementary Information*
232 (*Section A*).

233 3. Results

234 3.1. Sensor characterization tests: Example

235 Fig. 2 shows the data obtained in the first sensor characterization test
236 with sensor T1a on Oct. 6th, 2016 (day 3). The raw potential measurement
237 decreases during P0, increases to a steady value in P1, decreases to a steady
238 value in P2, increases to a steady value in P3, and decreases again in P4.
239 The time intervals used for computation of \tilde{E}^{P1} , \tilde{E}^0 , and \tilde{E}^{P3} (in calibration

240 medium, pH = 4, 7, and 4) are indicated by grey shading. One can see
241 that the measured offset \tilde{E}^0 is slightly below 0 mV (-1.30 mV). The values
242 for \tilde{E}^{P1} and \tilde{E}^{P3} are slightly lower than their ideal value (171.9 and 172.4
243 mV). The measured rise and decay sensitivities are therefore $\tilde{S}_D = 57.73$ and
244 $\tilde{S}_R = 57.90$ mV per pH unit. The results of every sensor characterization
245 test are visualized in the *Supplementary Information (Section G.2)*.

246 3.2. Long-term trends in the offset measurements within the warranty period

247 Fig. 3 displays the measured offsets in all sensors throughout the exper-
248 imental period. The recorded values collected within the warranty period
249 (1 year) range from approximately 0 mV (no offset) to roughly -70 mV.
250 All commercially available sensors (T1-T4) produce a decaying trend in the
251 offsets. The firstly recorded offsets for the T1-T3 sensors are small in magni-
252 tude and concentrate around 0 mV. In contrast, the T4 sensors offset values
253 indicate a shock effect producing a shift of -20 and -45 mV (T4a, T4b)
254 within days from installation. This is explained by the manufacturer as an
255 effect of the high ammonium concentration in the medium and should only
256 be expected for this specific type of sensors. The accumulated drift in the
257 T1 sensors is at most -25 mV after one year while the T2 and T3 sensors
258 exhibit an offset of -75 mV after one year. Without calibration, this means
259 the T1 sensors can produce a pH value as high as 7.4 when the true pH is 7.
260 The T2 and T3 sensors will produce a pH value as high as 8.3 in the same
261 circumstances. Due to failure of T5d, no offsets could be measured for this

262 sensor. The remaining prototypes (T5a/b/c) do not produce a significant
263 offset at any time, except for T5b which produces a dramatic shift in the
264 offset during three sensor characterization tests executed prior to replace-
265 ment. A detailed inspection of the T5b measurements revealed that the first
266 symptoms of sensor degradation can be observed on August 31st, 2017 (day
267 332). This is however only obvious when comparing these measurements
268 with the simultaneous T1b/T2b/T3b measurements (see the *Supplementary*
269 *Information, Section D*). In all cases, except for the T4 and T5a/b pairs, the
270 difference between offsets in sensors of the same type remains rather small
271 with 1 year of installation, with a maximal difference of 16.7 mV recorded
272 with the T2 sensors. Taking the 0.1 pH threshold discussed above as a
273 guideline, one could propose to validate and calibrate the sensors when their
274 potential measurements are 5.9 mV apart. This happens for the first time
275 for the T1, T2, and T3 sensors on day 127, 79, and 309. By these times,
276 the absolute offsets are already larger than this accepted threshold so that
277 the relative difference between sensors of the same type is unlikely a good
278 measure to trigger sensor maintenance.

279 Fig. 4 shows offsets for the sensors T1a, T3a, and T3b collected in the first
280 year of the experiment as a function of the difference in the offset between
281 T1a and T3a (left panel) and T3b and T3a (right panel). The left panel
282 suggests that offset difference between sensors can be predictive of the offset
283 in an individual sensor. The right panel shows that this is less likely to be
284 successful for sensors of the same sensor type, as also described above. This

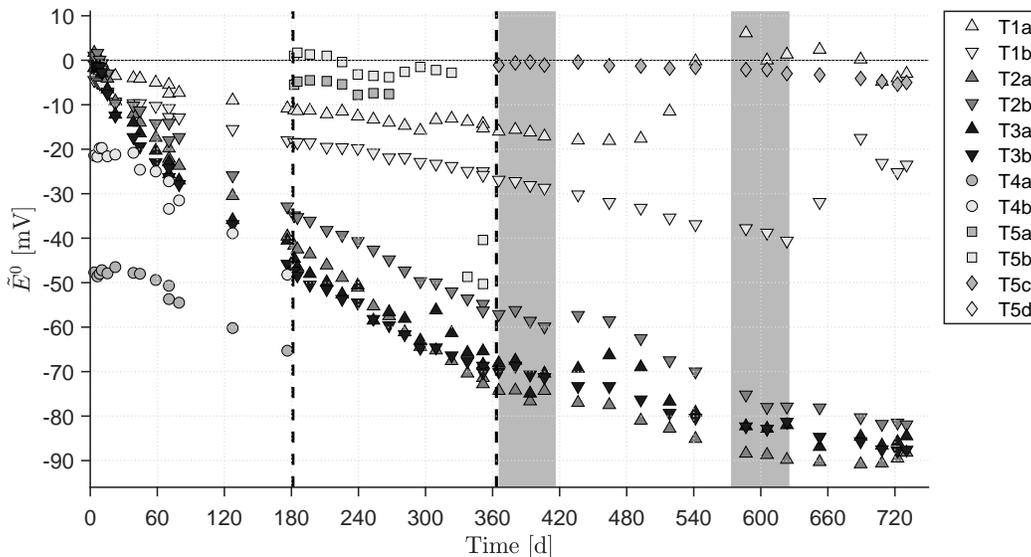


Figure 3: **Offset in all studied sensors as a function of time.** Vertical lines indicate a change of installed sensors (see Fig. 1). Grey bands indicate a change of reactor medium (see Section 2.3). The commercially available sensors (T1-T4) exhibit drift from the start of installation while the prototypes (T5) exhibit close to no drift when otherwise functioning properly. A significant shock effect is observed for the T4 sensors at the start of the experiment but not for any other sensor.

285 is considered an important opportunity for further research, which we discuss
 286 further below.

287 3.3. Long-term trends in the offset measurements beyond the warranty period

288 The offset measurements obtained after the warranty period expired ex-
 289 hibit two phenomena that are surprising (Fig. 3). The first phenomenon is
 290 the rise of the offset of the T1a sensor after 480 days of exposure and a similar
 291 rise of the offset of the T1b sensor after 630 days of exposure. Considering
 292 that this appears at distinct times in the lifetime of the T1 sensors, this can-
 293 not be explained as a direct effect of medium composition changes. Based on

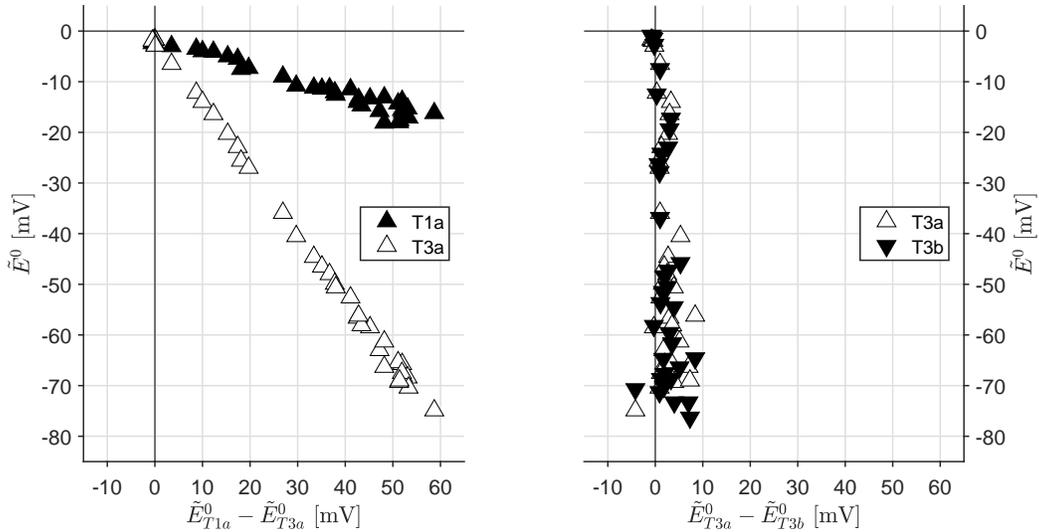


Figure 4: **Offset measurements as a function of relative deviations in the offset measurements.** *Left panel:* Offsets of sensor T1a and T3a as a function of the difference of these offsets. These data are suggestive of a close to linear relationship between sensor offsets and the offset difference. *Right panel:* Offsets of sensors T3a and T3b relative to the difference of these offsets. The difference in offset remains small and there is no obvious relationship in this case.

294 information provided by the sensor manufacturer, this type of drift rate sign
 295 reversal is unique for the T1 sensors and is unlikely to be observed with any
 296 other sensor type covered in this study. It is the opinion of the authors that
 297 the time for this reversal is difficult to predict in advance. For this reason,
 298 this phenomenon is best handled as an unmeasured process disturbance.

299 The second phenomenon consists of the rather flat to increasing profile of
 300 the offset measurements in the T2 and T3 sensors between day 360 and day
 301 480. Before and after this period, the drift rate in these sensors are visually
 302 similar. Given the synchronicity of this effect between 4 pH sensors, it is
 303 hypothesized that this change in the drift rate is influenced by the deliberate

304 addition of nitrite in the form of NaNO_2 salt to the reactor contents from day
305 366 to 417. The nitrite addition affected the biomass concentration and the
306 concentrations of all dominant nitrogen species (ammonia, nitrite, nitrate, see
307 *Supplementary Information, Section C*) and may also have affected the ion
308 strength and conductivity of the reactor contents. Due to this combination
309 of effects, the available data only offers an incomplete understanding of the
310 complete chain of causes and effects between the nitrite addition and the
311 observed changes in the sensor drift rates. For this reason, the effects of
312 changing media composition on the sensor drift rate is best also considered
313 an unmeasured process disturbance.

314 3.4. Long-term trends in the sensitivity measurements

315 Fig. 5 displays the computed sensitivity measurements for the potential
316 rise (\tilde{S}_R) during the complete experimental period. These measurements
317 do not exhibit strong trends in any particular direction. The sensitivity
318 measurements fall between 54.9 and 62.1 mV per pH unit. This means that
319 one can expect to measure a pH value between 5.95 and 6.08 when (i) the
320 true pH value is 6 and (ii) any offset is corrected for. The same graph also
321 shows the theoretical value of the sensitivity according to (2) and the recorded
322 temperature. This profile is very similar to the recorded sensitivity profiles
323 and explains most of the variations in the sensitivity measurements, which are
324 small anyway. The same conclusions are drawn from the computed sensitivity
325 measurements for the potential decay (\tilde{S}_D , see *Supplementary Information*,

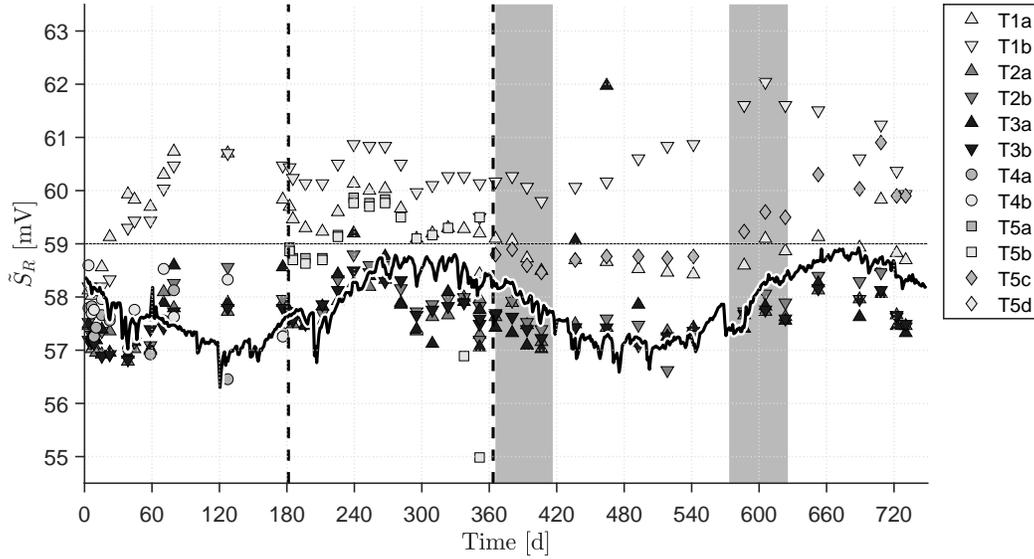


Figure 5: **Sensitivity measurements for the potential rise as a function of time.** Vertical lines indicate a change of installed sensors (see Fig. 1). Grey bands indicate a change of reactor medium (see Section 2.3). A black line shows the theoretically expected sensitivity computed with (2). Variations in the sensitivity are small and follow the theoretical sensitivity closely.

327 3.5. Drift models

328 For practical intents and purposes, the sensitivity – when corrected for
 329 temperature variations – can be considered constant for the considered pro-
 330 cess and sensors. We therefore focus on further analysis of the offset mea-
 331 surements.

332 The left panel of Fig. 6 shows the offset measurements for the T2a and
 333 T2b sensor together with the model predictions and their confidence bounds.
 334 The right panel of Fig. 6 shows the prediction residuals. With Model 1,

335 the time of the drift onset (t_f) is always identified as a time before the first
336 measurement was obtained (2.1 and 2.3 days), suggesting that drift occurs
337 throughout the experiment. The same kind of result is obtained with every
338 other commercially available sensor type (T1-T4), except for the T1a sensor
339 (see the *Supplementary Information (Section F)*). More importantly however
340 is that Model 1 offers a rather poor description of the data. The confidence
341 intervals are wide and the residuals are clearly auto-correlated. In contrast,
342 Models 2 and 3 provide narrower confidence intervals and residuals that do
343 not suggest presence of autocorrelation. There are no clear differences in
344 performance between these two models so that Model 3, which has fewer free
345 parameters, is preferred. The modeling results for the T1 and T3 sensors lead
346 to the same conclusions. For these results and all parameter estimates, we
347 refer to the *Supplementary Information (Section F)*. For the T4 sensors, all
348 model types delivered the same, adequate performance. This may indicate
349 that (a) the T4 sensors exhibit a drift which is influenced less by unmeasured
350 disturbances and therefore occurs with a close to constant rate or (b) that
351 the shortened exposure – 6 months in this case – was too short to capture
352 the long-term effects of unmeasured disturbances.

353 **4. Discussion**

354 This study present the first peer-reviewed results with which the effect
355 of long-term wear-and-tear on water quality sensors deployed in wastewater
356 treatment plants is assessed and evaluated in a systematic manner and at

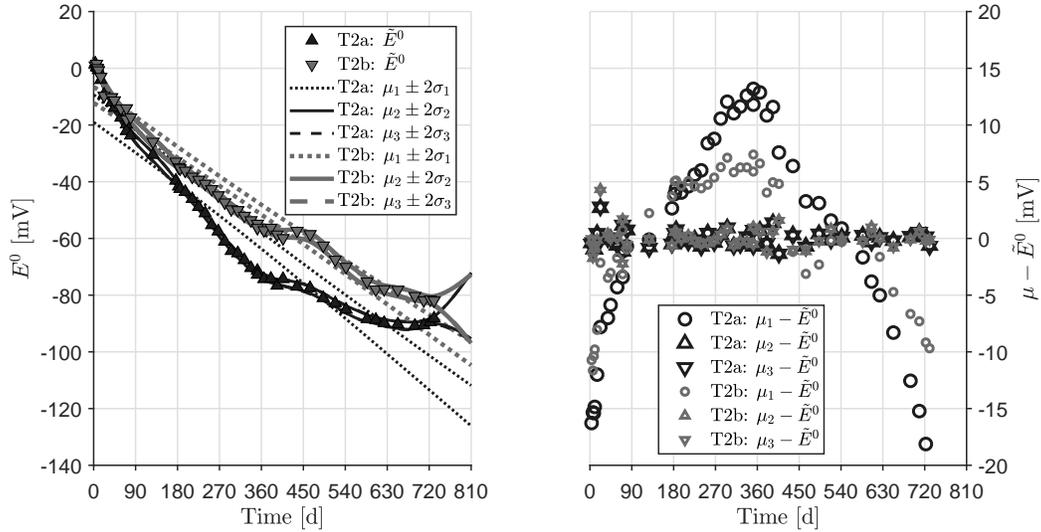


Figure 6: **Modeling results for the T2 sensors.** *Left panel:* Offset confidence bounds ($\mu \pm 2\sigma$) obtained with models 1 (μ_1, σ_1), 2 (μ_2, σ_2), and 3 (μ_3, σ_3). *Right panel:* Residuals between expected values (μ) and measured potentials (\tilde{E}^0). Model 1 does not describe the data well, leading to larger confidence bounds and auto-correlated residuals. Models 2 and 3 fit the data well and their predictions are hard to distinguish from each other.

357 this scale (12 sensors). The experimental results reveal that commonly held
 358 assumptions regarding the occurrence of sensors faults and fault symptoms
 359 are false. First, it is demonstrated that drift in pH sensors occurs simul-
 360 taneously in all commercially available sensors. Second, it is demonstrated
 361 that drift occurs as soon as a sensor is deployed in the measured medium.
 362 In some cases, the immediate onset of drift is paired by a significant shift in
 363 the offset. Importantly, the data needed to compute the offsets and sensi-
 364 tivities as a function of time are also available in modern pH instruments in
 365 the form of a calibration logbook that can be accessed through standardized
 366 communication protocols (e.g., Modbus).

367 These observations have important consequences for the development of
368 methods for fault detection and identification (FDI). Indeed, *(i)* one cannot
369 assume that faults appear independently in distinct sensors and *(ii)* one
370 cannot assume to have access to a fault-free historical data set. Naturally, this
371 also holds in the context of simulation-based benchmarking of FDI methods.
372 Consequently, it is our opinion that the development of FDI methods and
373 model-based benchmarking should be focused on methods that do not rely
374 on such assumptions.

375 Fortunately, our results also reveal a number of opportunities for the use
376 and maintenance of ion-selective measurements. First, the prototype sensors
377 tested in this study exhibit a remarkably stable offset. While these sensors
378 appear prone to failure, as one might expect from a prototype, this suggests
379 that practically drift-free yet economical pH sensors will enter the market
380 soon. Second, the recorded sensitivity measurements in all sensors hover
381 around the ideal values and are remarkably stable throughout the experimen-
382 tal period. Such a stable sensitivity lends support for advanced monitoring
383 and control strategies which are inherently robust to changes in the offset
384 but still assume a rather stable sensitivity [30, 24, 25]. Third, it was shown
385 that the offset difference between two pH sensors in the same medium can
386 be predictive of the offset of the individual pH sensors, however only if two
387 sufficiently distinct sensor types are selected. Combined with a stable sen-
388 sitivity, this means that the deviation between two online pH sensor signals
389 could be used as a proxy for the deviation in each individual sensor. Such

390 a proxy measurement could be very useful for remote sensor quality assess-
391 ment and predictive sensor maintenance, especially since one can compute
392 such deviations between on-line sensor signals while the sensors remain in
393 their normal measurement location in the monitored reactor.

394 The obtained offset measurements were studied in more detail by com-
395 paring the fit of 3 models. From this, it is concluded that the excessive
396 drift model included in the BSM family [18, 8] cannot adequately describe
397 the naturally occurring drift in ion-selective electrodes. Instead, the proposed
398 stochastic model, specifically an integrated Brownian process, delivers a good
399 description of the obtained data sets. In the authors' opinion, such a model
400 should be included in the BSM family for realistic simulation of measurements
401 obtained through ion-selective measurement principles. The obtained model
402 also enables prediction of the expected offset measurement and associated
403 confidence intervals beyond the last measurement. This means that such a
404 model can be used for predictive sensor maintenance, e.g., by planning a new
405 sensor validation and/or calibration before the predicted confidence interval
406 exceeds a predetermined tolerance, each time also updating the parameters
407 of the stochastic model. For this, confidence intervals for the reference po-
408 tential (E^0) rather than for the measurements (\tilde{E}^0) are expected to be most
409 useful. Exploring the utility of this idea is considered for future research.

410 **5. Conclusions**

411 Despite the abundance of literature of fault detection and identification
412 (FDI) methods, little is actually known about the cause-and-effect relation-
413 ships between the exposure of water quality sensors to harsh conditions, such
414 as wastewater media, and the occurrence of sensor faults and failures. This
415 first long-term study of the ageing of 12 individual pH sensors gives valuable
416 insight into this challenge. First, it is concluded that commonly held assump-
417 tions in FDI method development and evaluation, such as the availability of
418 fault-free historical data and independent onsets of sensor faults, are invalid
419 for pH sensors based on the ion-selective measurement principle. In addition,
420 the effects of offset drift in redundant sensors is unlikely to be identified early
421 if these sensors are of the exact same type and exposed to the same medium.
422 A stochastic model is shown to offer a good description of the observed drifts
423 of the sensor offsets and perform better than a previously established drift
424 model. Finally, our results suggest that newly developed pH sensors which
425 exhibit stable offsets will enter the commercial market soon.

426 **6. Acknowledgements**

427 This research was made possible by the Swiss National Foundation (Project:
428 157097). Mr. Ohmura's contributions are financially supported by Toshiba,
429 Tokyo, Japan. Dr. Carbajal's contributions are financially supported by
430 Eawag Discretionary Funds (grant no.: 5221.00492.011.05, project: DF2017/EMUmore).
431 All sensors were provided at no cost by Endress+Hauser. We thank Daniel

432 Iten and Stefan Vogel from Endress+Hauser for their valuable advice and
433 input to this study. The stochastic models were calibrated with the GPML
434 toolbox v4.2 [13].

435 **References**

- 436 [1] Alferes, J., Tik, S., Copp, J., Vanrolleghem, P. A., 2013. Advanced
437 monitoring of water systems using in situ measurement stations: data
438 validation and fault detection. *Water Science and Technology* 68 (5),
439 1022–1030.
- 440 [2] Barty, M., Patton, R., Syfert, M., de las Heras, S., Quevedo, J., 2006. In-
441 troduction to the DAMADICS actuator FDI benchmark study. *Control*
442 *Engineering Practice* 14 (6), 577–596.
- 443 [3] Buck, R., Rondinini, S., Covington, A., Baucke, F., Brett, C., Camoes,
444 M., Milton, M., Mussini, T., Naumann, R., Pratt, K., Spitzer, P., 2002.
445 Measurement of pH. definition, standards, and procedures (iupac rec-
446 ommendations 2002). *Pure and Applied Chemistry* 74 (11), 2169–2200.
- 447 [4] Carr, J. J., 1993. *Sensors and circuits*. PTR Prentice Hall.
- 448 [5] Corominas, L., Garrido-Baserba, M., Villez, K., Olsson, G., Cortes, U.,
449 Poch, M., 2018. Transforming data into knowledge for improved wastew-
450 ater treatment operation: A critical review of techniques. *Environmental*
451 *Modelling and Software* 106, 89–103.

- 452 [6] Corominas, L., Villez, K., Aguado, D., Rieger, L., Rosén, C., Vanrol-
453 leghem, P. A., 2011. Performance evaluation of fault detection methods
454 for wastewater treatment processes. *Biotechnology and Bioengineering*
455 108 (2), 333–344.
- 456 [7] Downs, J. J., Vogel, E. F., 1993. A plant-wide industrial process control
457 problem. *Computers and Chemical Engineering* 17 (3), 245–255.
- 458 [8] Gernaey, K. V., Jeppsson, U., Vanrollegheem, P. A., Copp, J. B., 2014.
459 Benchmarking of control strategies for wastewater treatment plants. *Sci-
460 entific and Technical Report No. 23. IWA Publishing.*
- 461 [9] Haimi, H., Mulas, M., Corona, F., Vahala, R., 2013. Data-derived soft-
462 sensors for biological wastewater treatment plants: An overview. *Envi-
463 ronmental Modelling and Software* 47, 88–107.
- 464 [10] Jeppsson, U., Pons, M. N., Nopens, I., Alex, J., Copp, J. B., Gernaey,
465 K. V., Rosén, C., Steyer, J., Vanrollegheem, P. A., 2007. Benchmark
466 Simulation Model No. 2: General protocol and exploratory case studies.
467 *Water Science and Technology* 56(8) (67-78).
- 468 [11] Le, Q. H., Verheijen, P. J., van Loosdrecht, M. C., Volcke, E. I., 2018.
469 Experimental design for evaluating WWTP data by linear mass bal-
470 ances. *Water Research* 142, 415–425.
- 471 [12] Nopens, I., Batstone, D. J., Copp, J. B., Jeppsson, U., Volcke, E., Alex,

- 472 J., Vanrolleghem, P. A., 2009. An ASM/ADM model interface for dy-
473 namic plant-wide simulation. *Water Research* 7, 1913–1923.
- 474 [13] Rasmussen, C. E., Nickisch, H., 2005. Gaussian processes for machine
475 learning (GPML) toolbox.
476 URL <https://gitlab.com/hnickisch/gpml-matlab/>
- 477 [14] Rasmussen, C. E., Williams, C. K., 2006. Gaussian processes for machine
478 learning. MIT press.
- 479 [15] Rieger, L., Langergraber, G., Siegrist, H., 2006. Uncertainties of spec-
480 tral in situ measurements in wastewater using different calibration ap-
481 proaches. *Wat. Sci. Technol.* 53(12), 187–197.
- 482 [16] Rieger, L., Takács, I., Villez, K., Siegrist, H., Lessard, P., Vanrolleghem,
483 P. A., Comeau, Y., 2010. Data reconciliation for wastewater treatment
484 plant simulation studies – planning for high-quality data and typical
485 sources of errors. *Water Environment Research* 82, 426–433.
- 486 [17] Rieger, L., Thomann, M., Gujer, W., Siegrist, H., 2005. Quantifying the
487 uncertainty of on-line sensors at WWTPs during field operation. *Water*
488 *Research* 39 (20), 5162–5174.
- 489 [18] Rosén, C., Rieger, L., Jeppsson, U., Vanrolleghem, P. A., 2008. Adding
490 realism to simulated sensors and actuators. *Water Science and Technol-*
491 *ogy* 57 (3), 337–344.

- 492 [19] Spindler, A., 2014. Structural redundancy of data from wastewater treat-
493 ment systems. determination of individual balance equations. *Water Re-*
494 *search* 57, 193–201.
- 495 [20] Spindler, A., Vanrolleghem, P. A., 2012. Dynamic mass balancing for
496 wastewater treatment data quality control using CUSUM charts. *Water*
497 *Science and Technology* 65 (12), 2148–2153.
- 498 [21] Taylor, B. N., Mohr, P. J., Douma, M., 2007. The NIST reference on
499 constants, units, and uncertainty.
500 URL physics.nist.gov/cuu/index
- 501 [22] Thomann, M., 2008. Quality evaluation methods for wastewater treat-
502 ment plant data. *Wat. Sci. Technol.* 10, 1601–1609.
- 503 [23] Thomann, M., Rieger, L., Frommhold, S., Siegrist, H., Gujer, W., 2002.
504 An efficient monitoring concept with control charts for on-line sensors.
505 *Water Science and Technology* 46 (4-5), 107–116.
- 506 [24] Thürlimann, C. M., Dürrenmatt, D. J., Villez, K., 2018. Soft-sensing
507 with qualitative trend analysis for control in full-scale wastewater treat-
508 ment plants. *Control Engineering Practice* 70, 121–133.
- 509 [25] Thürlimann, C. M., Udert, K. M., Morgenroth, E., Villez, K., 2018.
510 Assessment of two qualitative trend analysis tools for process control.
511 In: 4th IWA Specialized International Conference "Ecotechnologies for

- 512 Wastewater Treatment" (EcoSTP2018), London, ON, Canada, June 25-
513 27, 2018.
- 514 [26] Thürlimann, C. M., Udert, K. M., Morgenroth, E., Villez, K., Submit-
515 ted. Handling sensor drift for stabilizing nitrite control during nitrifica-
516 tion of high strength wastewater by means of qualitative trend analysis.
- 517 [27] Venkatasubramanian, V., Rengaswamy, R., Kavuri, S. N., 2003. A re-
518 view of process fault detection and diagnosis – Part II: Qualitative mod-
519 els and search strategies. *Computers & Chemical Engineering* 27 (3),
520 313–326.
- 521 [28] Venkatasubramanian, V., Rengaswamy, R., Kavuri, S. N., Yin, K., 2003.
522 A review of process fault detection and diagnosis – Part III: Process
523 history based methods. *Computers & Chemical Engineering* 27 (3), 327–
524 346.
- 525 [29] Venkatasubramanian, V., Rengaswamy, R., Yin, K., Kavuri, S. N., 2003.
526 A review of process fault detection and diagnosis – Part I: Quantitative
527 model-based methods. *Computers & Chemical Engineering* 27 (3), 293–
528 311.
- 529 [30] Villez, K., Habermacher, J., 2016. Shape anomaly detection for pro-
530 cess monitoring of a sequencing batch reactor. *Computers & Chemical*
531 *Engineering* 91, 365–379.

- 532 [31] Volcke, E. I., van Loosdrecht, M. C., Vanrolleghem, P. A., 2006.
533 Continuity-based model interfacing for plant-wide simulation: A gen-
534 eral approach. *Water Research* 15, 2817–2828.
- 535 [32] Westcott, C., 2012. pH measurements. Elsevier.
- 536 [33] Wolpert, D. H., 1996. The lack of a priori distinctions between learning
537 algorithms. *Neural Computation* 8 (7), 1341–1390.