

# Prediction of algal blooms via data-driven machine learning models: An evaluation using data from a well monitored mesotrophic lake

Shuqi Lin<sup>1,3\*</sup>, Donald C. Pierson<sup>1</sup>, Jorrit P. Mesman<sup>1,2</sup>

<sup>1</sup>Erken Laboratory and Limnology Department, Uppsala University, Uppsala, Sweden

<sup>2</sup>Département F.-A. Forel des sciences de l'environnement et de l'eau, Université de Genève, Genève, Switzerland

<sup>3</sup>Environment and Climate Change Canada, Canada Centre for Inland Waters, Burlington, ON, Canada, L7R 4A6

Correspondence to: Shuqi Lin (Shuqi.Lin@ec.gc.ca)

**Abstract.** With the increasing lake monitoring data, data-driven machine learning (ML) models might be able to capture the complex algal bloom dynamics that cannot be completely described in process-based (PB) models. We applied two ML models, Gradient Boost Regressor (GBR) and Long Short-Term Memory (LSTM) network, to predict algal blooms and seasonal changes in algal chlorophyll concentrations (*Chl*) in a mesotrophic lake. Three predictive workflows were tested, one based solely on available measurements, and the others applying a two-step approach, first estimating lake nutrients that have limited observations, and then predicting *Chl* using observed and pre-generated environmental factors. The third workflow was developed by using hydrodynamic data derived from a PB model as additional training features in the two-step ML approach. The performance of the ML models was superior to a PB model in predicting nutrients and *Chl*. The hybrid model further improved the prediction of the timing and magnitude of algal blooms. A data sparsity test based on shuffling the order of training and testing years showed the accuracy of ML models decreased with increasing sample interval, and model performance varied with training/testing year combinations.

## 1 Introduction

Harmful algal blooms, which are a serious threat to natural water systems, have been increasing throughout the world (Burford et al., 2020; Watson et al., 2016), primarily as a consequence of both climate change and increased nutrient loading from anthropogenic activities (Brookes and Carey, 2011; Paerl and Huisman, 2008). Moreover, as indicated by Carey et al. (2012) and Huisman et al. (2018), more intense and longer periods of thermal stratification could potentially specifically favour blooms of toxic cyanobacteria. To better manage and mitigate the effects of algal

26 blooms, methods to forecast their timing and magnitude are needed. However, the factors regulating algal blooms are  
27 complex, variable and site-specific, often involving high-order interactions of environmental factors and  
28 biogeochemical processes (Reichwaldt and Ghadouani, 2012; Richardson et al., 2018).

29 Process Based (PB) models encode our understanding of biogeochemical processes into a framework of numerical  
30 formulations, but these are inevitable simplifications that lead to an incomplete description of complex  
31 biogeochemical interactions and low level of model confidence (Elliott, 2012). Based on innovative data mining and  
32 statistical techniques, data-driven machine learning (ML) models have been applied to identify patterns within  
33 observed data (Peretyatko et al., 2012; Mellios et al., 2020), and with the recent proliferation of lake monitoring data  
34 (Marcé et al., 2016), ML models have been applied, as an alternative to PB models for bloom prediction (Rousso et  
35 al., 2020). Previously applied ML models, including Random Forest (Nelson et al., 2018), Support Vector Machine  
36 (Jimeno-Sáez et al., 2020), and Artificial Neural Network (Xiao et al., 2017; Recknagel et al., 1998; Wei et al., 2001),  
37 can improve predictions of the timing and seasonality of algal *Chl* pattern, apparently by accounting for complexity  
38 that is difficult to encode within the framework of a PB model. However, a downside of data-driven ML models is  
39 that they lack the interpretability and generalization found in the explicit structure of the PB model. In recent years,  
40 process-guided-deep learning (PGDL) model emerged and was applied to water temperature (Jia et al., 2019; Read et  
41 al., 2019) and water quality (Hanson et al., 2020) simulations, which explicitly combine well-defined physical theories  
42 into the training of ML models, enhancing their interpretability. While this approach has achieved promising results,  
43 it is difficult to apply it to phytoplankton dynamics due to numerous nonlinear interactions within the biogeochemical  
44 cycles and the difficulty in defining a measurable processes or mass balances that can be used as a physical constraint  
45 on knowledge-guided decisions. Also, the sparsity of lake water quality (e.g., nutrients, *Chl* concentration)  
46 observations can limit the application of ML models in algal bloom modelling (Rousso et al., 2020).

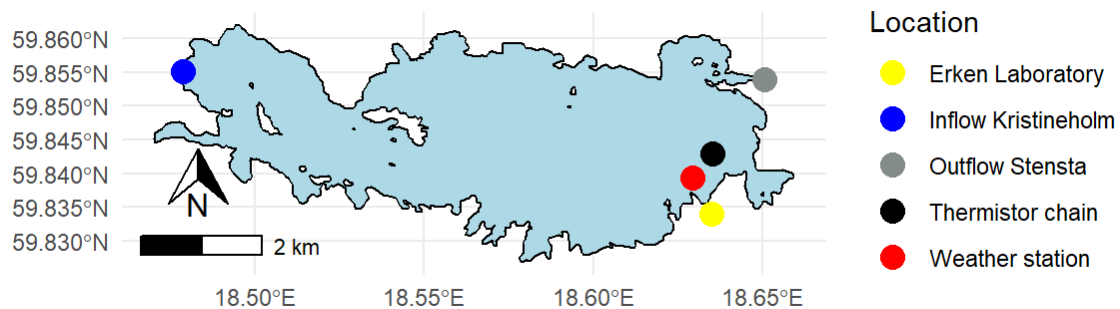
47 In this study, our objectives are to (1) apply the ML models to predict algal bloom in a well-monitored mesotrophic  
48 lake; (2) evaluate model performance and assess model uncertainties; (3) explore the approaches to improve the model  
49 performance and widen the model applications. We first tested the ability of ML models in predicting algal *Chl*  
50 concentrations via available environmental factors, including observed lake nutrients data, and then proposed a two-  
51 step ML approach for predicting algal dynamics that: first estimates lake nutrient concentrations which often have  
52 limited observations and secondly predicts variations in algal *Chl* using these pre-generated nutrient concentrations  
53 combined with other observed environmental factors that are collected at higher frequency. We also tested a simple

54 hybrid model architecture that by adding hydrodynamic features derived from the PB model into the training features  
55 of the two-step ML approach, allowing us to include additional information describing physical lake processes  
56 expected to affect variations in algal growth and succession in the machine learning prediction.  
57 We applied the above workflows to predict changing *Chl* concentration, as a proxy for the occurrence of algal blooms,  
58 via Gradient Boost Regressor (GBR) and Long Short-term Memory network (LSTM). Two shuffling year tests were  
59 conducted. One assessed the uncertainty of ML models in predicting *Chl* during the same two-year period and the  
60 other evaluated the sensitivity of ML accuracy to various training/testing year combinations and lake nutrient sampling  
61 intervals. Model performance and potential applications in algal bloom forecasting are discussed.

## 62 2 Methods

### 63 2.1 Study site

64 The study site, Lake Erken, is a mesotrophic lake located in east-central Sweden, that has a surface area of 24 km<sup>2</sup>, a  
65 maximum depth of 21 m and an average retention time of 7 years. The lake is dimictic with seasonal stratification  
66 commonly beginning in May-June and ending in August-September. The onset of ice cover usually begins in  
67 December-February and the loss of ice occurs in Mar-April (Persson and Jones, 2008). Located near the Baltic coast,  
68 Lake Erken is wind exposed, and susceptible to periodic wind-induced turbulent mixing.  
69 Changes in algal *Chl* in Lake Erken have a typical seasonal pattern, with spring and summer peaks in concentration  
70 (Pettersson et al., 2003). Spring blooms are dominated by dinoflagellates and diatoms (Pettersson, 1985), and initiated  
71 by overwinter species from the last autumn (Yang et al., 2016). Cyanobacteria dominate summer peaks in *Chl*, given  
72 that they can optimize their vertical position in regarding to nutrients and light (Paerl, 1988; Pierson et al., 1992).



73

74 **Figure 1.** Map of Lake Erken. The locations of the monitoring systems are shown.

## 75 2.2 Data

76 Lake Erken has a long running automated monitoring program that provides hourly meteorological data, water  
77 temperature profiles between 0.5 and 15 m at 0.5 m intervals and the flow from the inflow and outflow (Fig.1). A  
78 manual sampling program collects samples during ice-free time at 5-7 days intervals for all major nutrient  
79 concentrations (e.g., NO<sub>x</sub>, NH<sub>4</sub>, PO<sub>4</sub>, Total P, Si, etc.), dissolved oxygen (O<sub>2</sub>), and *Chl* concentration. The timing of  
80 the onset and loss of ice cover are also monitored yearly by the lab. More detailed information on the sampling program  
81 is in Supporting Information (See Text S1) and Moras et al. (2019).

## 82 2.3 Modelling Methods

### 83 2.3.1 Process-based (PB) lake model

84 In this study, a PB hydrodynamic lake model, GOTM (General Ocean Turbulence Model) (Burchard et al., 1999),  
85 was used to generate water temperature profiles, and other hydrodynamic metrics. GOTM also served as the  
86 foundation of water quality simulations made with the SELMAPROTBAS model (Mesman et al., 2022) that is coupled  
87 to GOTM through the Framework for Aquatic Biogeochemical Models FABM (Bruggeman and Bolding, 2014).

### 88 2.3.2 Data-driven machine learning (ML) models

89 Tree models have been widely applied in modelling phytoplankton dynamics in freshwater systems (Harris and  
90 Graham, 2017; Fornarelli et al., 2013; Rousso et al., 2020). Gradient Boosting Regressor (GBR) is one of these tree  
91 models, iteratively generating an ensemble of estimator trees with each tree improving upon the performance of the  
92 previous. The details about GBR model can be found in Friedman (2001). The hyperparameters in GBR are optimized  
93 via *RandomizedSearchCV* function within Scikit-Learn library. The loss function of model is chosen as ‘huber’, which  
94 is a combination of the squared error and absolute error of regression. Since the target variable in our research *Chl*  
95 concentration has peak values during algal blooms which could be regarded as outliers, the ‘huber’ loss function is  
96 more robust and gives greater weight to peak values than the mean squared error function.

97 Long short-term memory (LSTM) network is part of a class of deep learning architectures, called recurrent neural  
98 network (RNN), built for sequential and timeseries modelling (Hochreiter and Schmidhuber, 1997). The core concepts  
99 of LSTM are the cell and hidden states, and its three gates (input gate, forget gate, and output gate; See Fig. S2).  
100 Essentially, the LSTM model defines a transition relationship for a hidden representation through a LSTM cell which  
101 combines the input features at each time step with the inherited information from previous time steps. This architecture

102 is suitable for extracting information from sequential data (Rahmani et al., 2020; Read et al., 2019). The  
103 hyperparameter settings in LSTM can be found in Supporting Information (See Text S2).  
104 Compared to GBR model, LSTM has more complex model architectures, carrying the ‘memory’ from the previous  
105 time steps. In this study, GBR and LSTM were applied, respectively, to assess the performance of ML models with  
106 and without ‘memory’. Both ML models are built in Python using the Scikit-Learn (<https://scikit-learn.org/stable/>, last  
107 access: September, 2022) and TensorFlow (<https://www.tensorflow.org/>, last access: September, 2022) libraries.

#### 108 **2.4 Design of predictive workflows and shuffling year data sparsity tests**

109 In this study, we tested three workflows using a dataset split for training (years 2004-2016) and testing (years 2017-  
110 2020). In all three workflows, a 5-fold cross-validation using the training dataset was used to optimize the  
111 hyperparameters in the ML models. Workflow 1 directly predicts *Chl* concentration based on available environmental  
112 observations (Table 1). The training and testing datasets were limited by the frequency of lake nutrient observations  
113 which resulted in 5-7 day gaps between data points. The time step of LSTM was set to 1, that is, the environmental  
114 factors on the target date and previous observation date, which may be 5-7 days ago, were used to train the model and  
115 make predictions.

116 In workflow 2 and 3, a two-step approach was applied (Table 1). Daily measurements of physical factors were used  
117 to pre-generate daily variations in lake nutrients via separate ML models, and the ML models were trained at a daily  
118 time step using the measured environmental factors and pre-generated nutrient concentrations. The time step of LSTM  
119 was then set to 7 days.

120 In workflow 3, three hydrodynamic features, i.e., mixing layer depth ( $z_e$ ), Wedderburn number ( $W_n$ ), and the seasonal  
121 thermocline depth (*thermD*), derived from the GOTM model were regarded as daily training features in the two-step  
122 ML approach. The definitions and calculations of these features are explained in SI (2.5 Feature selection and  
123 processing for ML models, Text S3)

124 Following the two-step approach and using workflow 3, we set up two tests. (1) To assess the uncertainty induced by  
125 variations in the data used to train the ML models, we shuffled the training years, randomly taking 13 years out of  
126 2004-2018 dataset 30 times, and tested the model predictions of *Chl* during 2019-2020. And, (2) to test if the workflow  
127 could be used for other water systems which may have less frequent lake nutrient monitoring data, we conducted a  
128 data sparsity test that evaluated the sensitivity of models to the lake nutrient and *Chl* sampling interval. For this test  
129 the lake nutrient and *Chl* concentration observations in training dataset was down-sampled to a 7-day, 14-day, 21-

130 day, 28-day, and 35-day sampling interval. Then for each sampling interval using the 2004-2020 dataset, *Chl* was  
 131 predicted for different consecutive 4-year periods when the ML models were trained by the remaining 13 years of  
 132 data. Data shuffling was conducted 13 times so that every 4-year period in our dataset was tested.

133 **Table 1** List of training features and target variables in each workflow. Blue indicates training features, red indicates  
 134 target variables, purple indicates the variables are the target variables in step 1 used to produce daily a training feature  
 135 for use in step 2. The order of nutrient model sequence is from the top to bottom based on its position in the table  
 136 (NO<sub>x</sub> to Si).

variables	Sample interval	workflow 1	workflow 2		workflow 3	
			Step 1	Step 2	Step 1	Step 2
Inflow	Daily	[Blue shaded area]	[Blue shaded area]	[Blue shaded area]	[Blue shaded area]	[Blue shaded area]
Meteorological data (Air temperature, wind speed, shortwave radiation, precipitation, humidity, cloud cover)	Daily					
$\Delta T$	Daily					
Ice duration	Daily					
Days from ice-off date	Daily					
$z_e$	Daily					
$W_n$	Daily					
<i>thermD</i>	Daily	[Blue shaded area]	[Blue shaded area]	[Blue shaded area]	[Blue shaded area]	[Blue shaded area]
NO <sub>x</sub>	1-2 weeks					
O <sub>2</sub>	1-2 weeks					
PO <sub>4</sub>	1-2 weeks					
Total P	1-2 weeks					
NH <sub>4</sub>	1-2 weeks					
Si	1-2 weeks					
Chl	1-2 weeks	[Red shaded area]	[Red shaded area]	[Red shaded area]	[Red shaded area]	[Red shaded area]

137

### 138 2.5 Feature selection and processing for ML models

139 The feature selection process is based on some a priori knowledge of the underlying phenomena related to algal  
 140 blooms. All workflows made use of the daily automated monitoring data. In addition, the temperature difference ( $\Delta T$ )  
 141 between surface water (averaged over the upper 3 m) and bottom water (15 m) was also used to represent the thermal  
 142 structure of the lake., and the duration of ice cover in the previous winter, and the number of days from ice-off date  
 143 were used.

144 In workflow 2 and 3 nutrients are predicted sequentially, with each pre-generated nutrient predictions included in the  
 145 training data of the next nutrient prediction (Table 1). Workflow 3 added  $z_e$ , computed using the GOTM simulated  
 146 vertical eddy diffusivity ( $K_z$ ) profiles, *thermD*, estimated using Lake Analyzer (Read et al., 2011) based on GOTM  
 147 simulated temperature profile, and  $W_n$ , a dimensionless parameter measuring the balance between wind stress and the  
 148 pressure gradient resulting from the slope of the interface (See Text S3, SI), as additional daily training features.

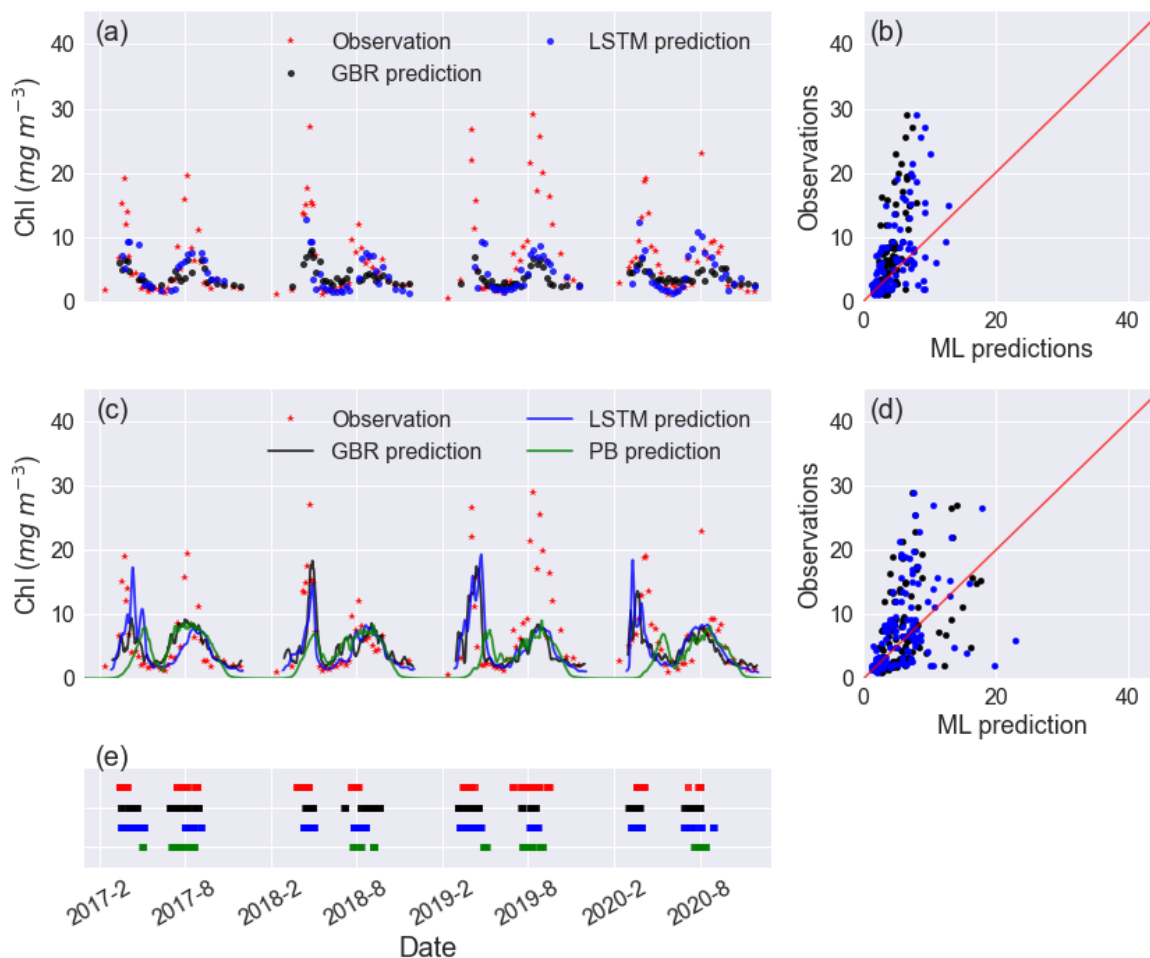
## 149 **2.6 Evaluating metrics**

150 Model performance was evaluated by comparing the simulated and measured *Chl* concentrations, and by calculating  
151 the mean absolute error (*MAE*), root means square error (*RMSE*), and correlation coefficient ( $R^2$ ). To evaluate the  
152 accuracy of the model in detecting the onset of an algal bloom, we calculated a confusion matrix in workflows 2 and  
153 3, where the observations were linearly interpolated to daily values, and predicted daily *Chl* concentration were  
154 smoothed with a 7-day rolling mean. Using these data, the onset of a bloom was categorized as occurring when the  
155 daily change of *Chl* ( $\Delta Chl$ ) exceeded a threshold,  $0.35 \text{ mg m}^{-3} \text{ day}^{-1}$ . This works well in Lake Erken where *Chl*  
156 concentrations are frequently monitored (near weekly), and the linear interpolation can be expected to be reasonably  
157 representative of the *Chl* concentrations between measured samples. Considering the randomization in the ML models,  
158 we also add a 3-day window on the bloom onset prediction, that is, we considered the prediction of a bloom valid if  
159 the measured data suggested a bloom the day before or after the simulated onset. We used the True Positive Rate  
160 (TPR), False Positive Rate (FPR), and modified accuracy (Kappa) which considers the possibility of the agreement  
161 occurring by chance (McHugh, 2012), to identify the potential of ML models to correctly capture the algal bloom  
162 onset (See Table S1, SI). A model with 100% TPR, 0% FPR, and 100% Kappa would constitute a perfect fit.

## 163 **3 Results**

### 164 **3.1 Workflow 1: Direct prediction based on observations**

165 In workflow 1, both GBR and LSTM clearly reproduced spring and summer blooms (Fig. 2a) but underestimated the  
166 intensity of blooms (Fig. 2a, b). Neither ML model captured the extraordinarily high *Chl* ( $\sim 15\text{-}30 \text{ mg m}^{-3}$ ) in the  
167 summer of 2019. Although the abnormal summer bloom in 2019 could contribute to the higher *RMSE* and *MAE* in the  
168 testing dataset than the mean values in the training dataset, the cross-validation on the training dataset (See Table S2,  
169 SI) shows what appears possibly to be overfitting issue in both models. The achieved accuracy of models is attributed  
170 to the daily availability of physical inputs, and the fact that in Lake Erken water samples are collected frequently at 5-  
171 7 days intervals. Workflow 1 may be most valuable in reconstructing previous variations in algal *Chl*, filling the gaps  
172 between measured *Chl* observations and feature importance ranking (See Fig. S4, SI). But when using this workflow,  
173 future forecasts will be limited by the absence of future nutrient data.



174  
 175 **Figure 2.** Timeseries of observed and predicted *Chl* from GBR and LSTM models in (a) workflow 1 and (c) workflow  
 176 3, and the corresponding scatter plots of observations vs ML predictions of *Chl* in workflow 1 and workflow 3 are  
 177 shown in panels (b) and (d), with the black and blue dots/lines representing the predictions from GBR and LSTM,  
 178 respectively. Panel (e) shows the observed and predicted algal bloom onsets in 2017-2020 using the same color coding  
 179 as the previous panels. Results from the PB model simulation in Mesman et al. (2022) are also shown in (c) and (e).

### 180 3.2 Workflow 2: Two-step ML models based on pre-generated daily nutrients and observed physical factors

181 As in workflow 1, both ML models in workflow 2 had poor fit in the summer of 2019 and suffered from overfitting  
 182 leading to higher *MAE*, *RMSE*, and lower *R*<sup>2</sup> in testing datasets than training datasets (See SI, Table S2).

183 Overall, both GBR and LSTM showed slightly higher *MAE* (4.22 mg m<sup>-3</sup> vs. 3.87 mg m<sup>-3</sup>) and *RMSE* (6.27 mg m<sup>-3</sup>  
 184 vs. 6.00 mg m<sup>-3</sup>) when compared to workflow 1 (Table 2). But they also showed improved performance in terms of  
 185 capturing the peak values of *Chl* during spring blooms (Fig. 2, Fig. S5, SI). Both workflows outperformed the  
 186 SELMAPROTBAS PB model in simulating concentrations of lake nutrients (See Fig. S6, SI). The ML models were



187 more accurate in predicting the low values of NO<sub>x</sub> and peak values of PO<sub>4</sub> and Total P. However, both ML models  
 188 and the PB model failed in predicting the extremely high values of measured lake nutrients, such as the autumn peak  
 189 of NH<sub>4</sub> in 2017 (Fig. S6e) and the spring peak of O<sub>2</sub> in 2018 (Fig. S6c), Thus, higher workflow 2 *MAE* and *RMSE*  
 190 (Table 2) are presumably due to the inaccuracies in the pre-generated nutrient training data, but the improved daily  
 191 predictions that better capture the bloom events, overshadow these flaws.

192 **Table 2** Comparisons of model performance during the testing period based on *RMSE*, *MAE*, and *R2*. The unit of *Chl*  
 193 is mg m<sup>-3</sup>. In bold are the best fits of each statistical metric. For comparison of training and testing periods, see Table  
 194 S2.

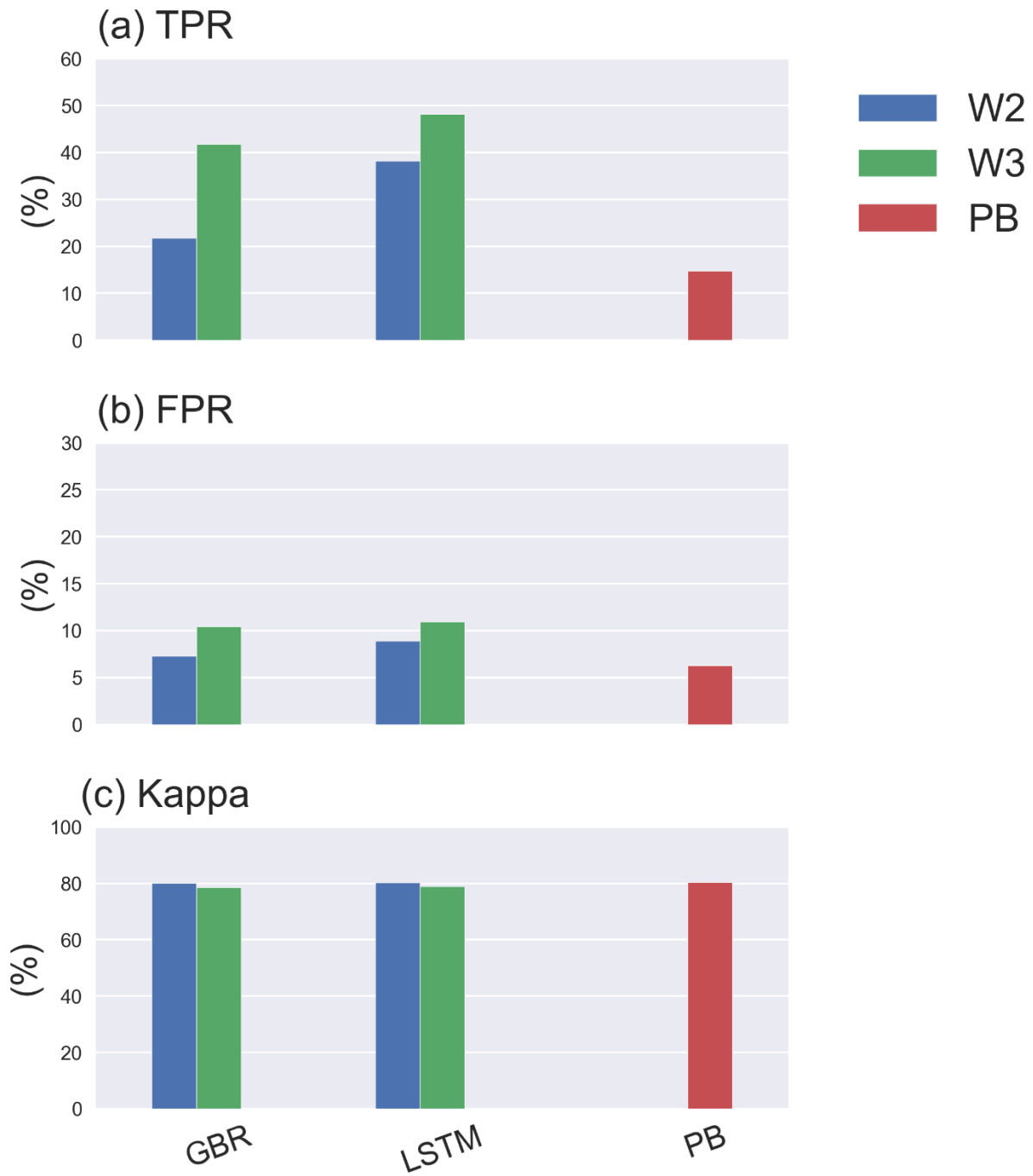
Model	PB	ML-workflow 1		ML-workflow 2		ML-workflow 3	
		GBR	LSTM	GBR	LSTM	GBR	LSTM
<i>RMSE</i>	7.18	5.77	<b>5.64</b>	6.27	6.00	5.94	5.81
<i>MAE</i>	4.77	<b>3.55</b>	3.58	4.22	3.87	3.99	3.71
<i>R2</i>	-0.25	0.13	<b>0.20</b>	0.05	0.13	0.14	0.18

195  
 196 **3.3 Workflow 3: based on workflow 2, and including hydrodynamic training features derived from the**  
 197 **GOTM model.**

198 Including hydrodynamic training information in workflow 3 did not significantly improve in lake nutrient predictions  
 199 compared to workflow 2 (See Fig. S6), and when using workflow 3 both ML models showed comparable performance  
 200 in *Chl* predictions compared to workflow 1. However, the predictions of the spring bloom in all years improved  
 201 compared to workflows 1 and 2, in terms of the magnitude and timing of the spring bloom (Fig. 2e). This was the case  
 202 in 2019-2020 (Fig. 2a) which was an abnormally warm winter with only 5 days ice cover, and had an unusually early  
 203 spring algal bloom. Both GBR and LSTM in workflows 2 and 3 did not capture the extremely intensive bloom (with  
 204 peak values close to 30 mg m<sup>-3</sup>) in summer of 2019, and neither did the PB model.

205 Furthermore, adding hydrodynamic features derived from PB model improved predictions of the onset of algal blooms  
 206 (Fig. 2e and 4), with the overall TPR increasing by 15 % and 5 %, FPR increasing around 5% and 3 % in GBR and  
 207 LSTM models, respectively. Compared with the PB model which showed lower TPR (15%) and FPR (6%), ML  
 208 models are more likely to predict algal bloom at the correct time. The optimal TPR was from LSTM in workflow 3,  
 209 which could detect the onset of algal blooms with TPR closed to 50%. However, the concomitant higher FPRs  
 210 indicating an incorrect warning of algal bloom is also more likely to occur in the ML models, since the PB model is  
 211 more like to miss the bloom entirely. The Kappa values of both ML models and the PB model are close to 80%,

212 showing that all models simulated the entire period (blooms and the periods between blooms) to a moderate-strong  
213 level (McHugh, 2012).



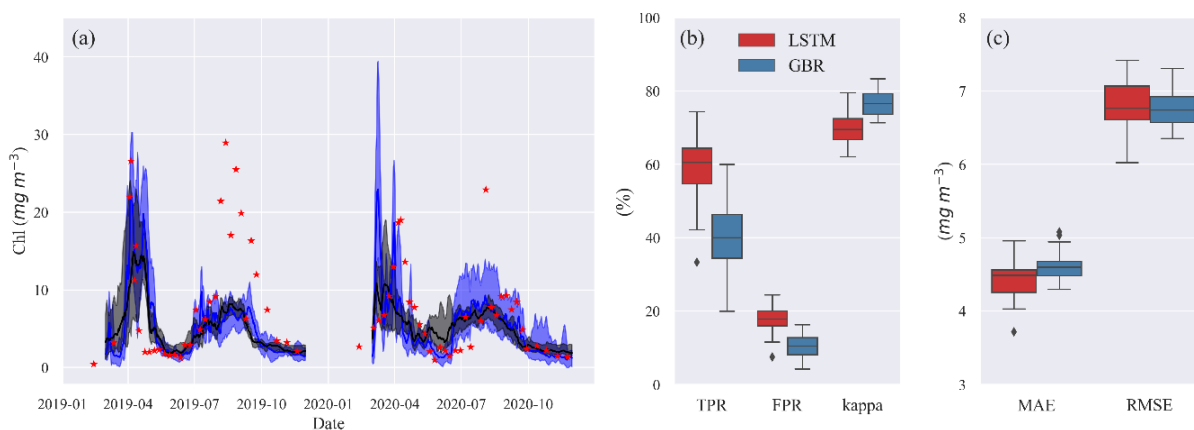
214

215 **Figure 3.** TPR, FPR, Kappa of GBR and LSTM models in workflow 2, 3 and the PB model.

### 216 3.4 Effects of shuffling training years on 2019-2020 predictions

217 The results presented so far are based on a typical strategy of training ML models for a historical period in this case  
218 2004-2016 and then accessing model performance in a second period between 2017-2020. The accuracies of the model  
219 predictions were to some extent related to the range and variability in the training data. To evaluate the importance of  
220 this we randomly removed two years from a 2004-2018 training dataset, and made 30 different predictions of *Chl*  
221 during 2019-2020 when the models had difficulties predicting spring and summer blooms (Fig 5). When trained with  
222 the various shuffled combinations, both ML models were capable of reproducing the seasonal variations in algal *Chl*  
223 with a 4.5 % and 5.8 % coefficient of variation (CV) in MAE, and a 24.0 % and 16.4 % CV in TPR of GBR and LSTM,  
224 respectively (See Table S3, SI). This provides an indication of the uncertainty that may arise as a consequence of  
225 differences in the training datasets used for in our workflows. And, it also shows that even a relatively long training  
226 period of 13 years can not totally capture the system behaviour in such a way as to lead to nearly similar bloom  
227 predictions.

228 Although none of the model runs captured the intensive summer bloom in 2019, the spring bloom in both years was  
229 well represented, especially by LSTM, in terms of timing and magnitude.



230 **Figure 4.** (a) Timeseries of observed (red stars) and predicted *Chl* from GBR (black) and LSTM (blue) models in the  
231 shuffling training year test. The shades represent the range between minimum and maximum prediction, and the solid  
232 lines represent the median prediction. (b) shows the boxplot of TPR, FRP, and Kappa, and (c) shows boxplot of MAE  
233 and RMSE of both models in the shuffling training year test.

234  
235  
236 Despite comparable *RMSE* and *MAE* in LSTM and GBR (Fig. 4c), both higher TPRs (with median of 60%) and FRPs  
237 (with median of 18%) in LSTM indicate that the LSTM was more aggressive in making algal bloom predictions. The  
238 GBR model's apparent advantage in FPRs (with median 10%) is largely the result of it making a lower number of

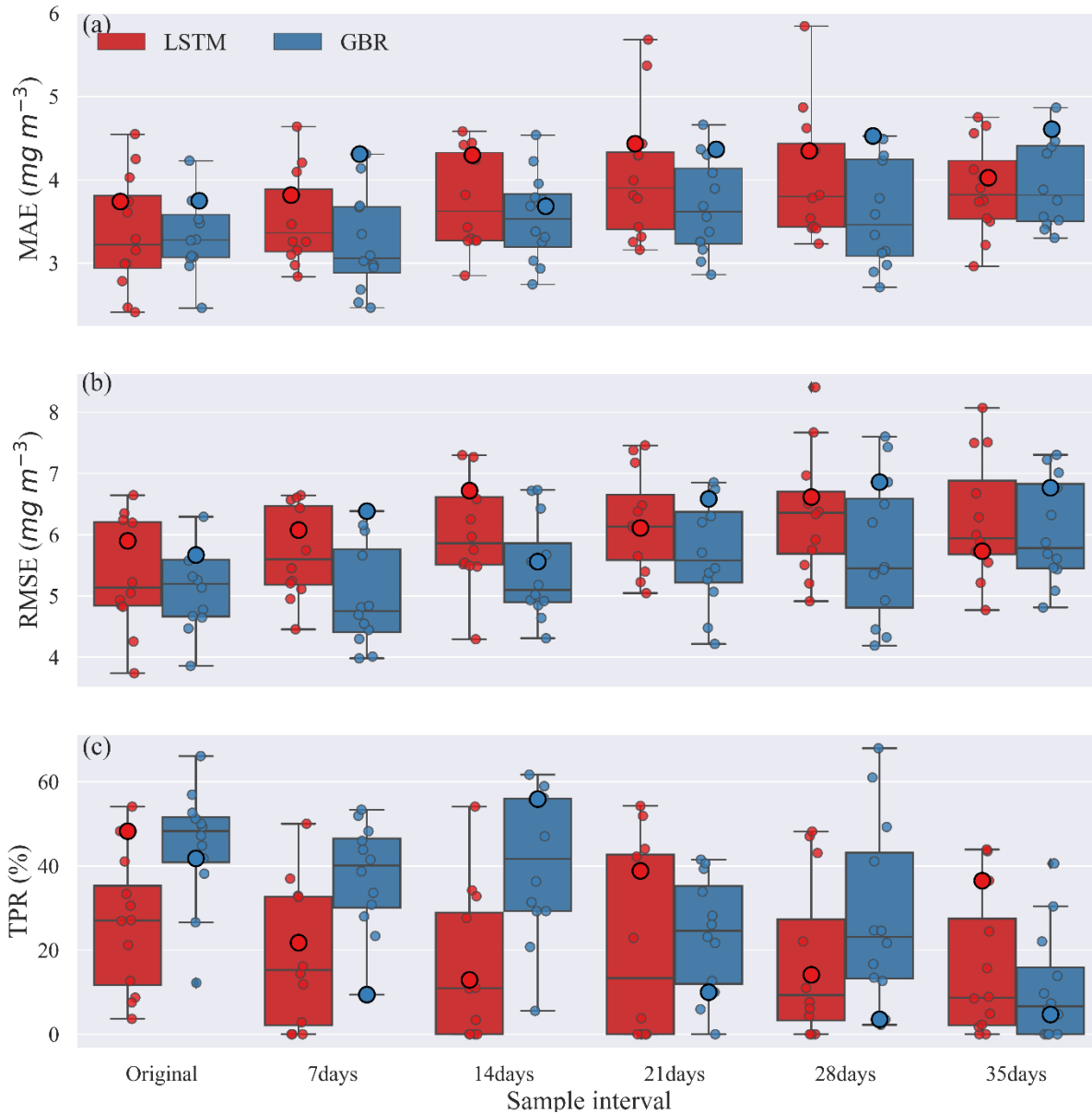
239 bloom predictions since the low concentrations between spring and summer blooms in 2020 was not well represented  
240 (Fig. 4b).

### 241 **3.5 Shuffling years data sparsity test**

242 To examine the possible use of workflow 3 when data are less frequently available, lake nutrient and *Chl* data were  
243 down-sampled so that the effects of sampling frequency on model predictions could be evaluated. Each down-sampled  
244 dataset was also rearranged into 13 different 13-year training periods and 4-year testing periods. The variability in  
245 predictions provided a measure of model performance and uncertainty. Fig. 5 shows the uncertainty in model  
246 predictions as a consequence of the chosen sampling intervals.

247 The *MAEs* and *RMSEs* of both GBR and LSTM models tended to increase with the longer sample intervals. The  
248 median *MAE* was always slightly higher for the LSTM model except when trained with original dataset (Fig. 5a).  
249 While our initial evaluation of TPR using 2017-2020 as the testing period and 2004-2016 as the training period  
250 suggested the LSTM model was more accurate in turns of detection of algal bloom onsets (Fig. 3), Fig. 5c showed the  
251 median TPR of GBR model calculated by the shuffling year test was over 50%, higher than that found when using the  
252 original testing and training periods. This can be explained by the fact that the 2017-2020 testing period as in Fig. 3  
253 and shown as large points in Fig. 5 was unusually difficult for GBR to simulate. Consequently, even though the GBR  
254 model usually performs better in the shuffled data test in Fig. 5, Fig. 3, which shows the results of 2017-2020 testing  
255 period, presented the opposite result. This illustrates the importance of the sequence of training and testing years for  
256 evaluating model performance.

257 For the first three sampling intervals the GBR model clearly had better TPR values than the LSTM model. The median  
258 TPRs of GBR model started to drop below 30% once the sample interval reached 21 days. For LSTM, median TPRs  
259 remained lower than 30%, for all sampling intervals but also showed a much wider range of variability (Table S4)  
260 dependent on the training and tested datasets used. In general, both models performed best at the original and 7-day  
261 sampling interval, but then showed slightly worse performance that was consistent up to a sample interval of 21 days.  
262 In terms of the errors evaluated over the entire 4-year testing period (Fig. 5a, b) the GBR model had lower errors and  
263 therefore, better predicted the seasonal variations of *Chl* concentration. The timeseries comparison of observed and  
264 predicted *Chl* from this shuffling year data sparsity test can be found in SI (Fig. S7-9).



265

266 **Figure 5.** Comparisons of (a) *MAE*, (b) *RMSE*, and (c) *TPR* between GBR and LSTM during the testing period created  
 267 under various sample intervals. Circles along the box show the result from the testing period of all shuffled  
 268 training/testing year combinations and the bigger circles represent 2004-2016 training and 2017-2020 testing years  
 269 combination as was used in Fig. 2.

## 270 4 Discussion

### 271 4.1 Performance of ML models

272 In three workflows, the ML models successfully reproduced the *Chl* seasonal patterns, capturing the spring and  
 273 summer bloom events, with lower averaged *RMSEs* and *MAEs* than a PB model simulation that was previously

274 calibrated for Lake Erken. And in all three workflows, LSTM model always showed slightly lower *RMSE*, *MAE* and  
275 higher *R*<sup>2</sup> in predicting *Chl* concentrations than GBR model, and higher TPR in detecting the onset of algal bloom  
276 events. Workflow 1 which predicted *Chl* based on all available environmental factors including lake nutrient  
277 observations showed that both ML models can reproduce the seasonal dynamics of algal *Chl* with promising accuracy  
278 (*MAE* = 3.55 and 3.58 mg m<sup>-3</sup>, *RMSE* = 5.77 and 5.64 mg m<sup>-3</sup> and *R*<sup>2</sup> = 0.13 and 0.20, for GBR and LSTM, respectively)  
279 via the direct input of available environmental observations. These ML models can be applied to reconstruct past  
280 patterns of algal *Chl*, fill the gaps between measured *Chl* observations, and interpret the mechanisms that drive  
281 phytoplankton dynamics. Workflows 2 and 3 adopted a two-step approach, first using separate ML models to  
282 estimating daily changes in lake nutrient concentration, and in Workflow 3 also including PB model derived physical  
283 factors as training features of the algal ML model. These two workflows allowed daily predictions of changes in algal  
284 *Chl* concentration using both observations and pre-generated lake nutrient concentrations at a consistent daily time  
285 step, and at only a minor decrease in performance compared to workflow 1, workflow 2 and 3 demonstrated a wider  
286 potential range of applications (e.g., interpolation, reconstruct historical data, algal bloom forecast) via making daily  
287 forecasts with less-than-daily measured nutrient observations.

288 The one clear failure of both the ML and PB based model predictions was during July-August 2019, *Chl* concentrations  
289 in integrated samples collected between the surface and 6-12 m exceeded 20 mg m<sup>-3</sup> over a 5-week period. Neither  
290 the PB model nor ML models captured this unusually persistent bloom (Fig. 2, Fig. S3, SI). At this time the  
291 phytoplankton were dominated by the cyanobacteria *Gloeotrichia* and *Anabaena*, that form a resting akinete life stage  
292 at the end of their yearly bloom, which can initiate the following year's bloom as they are transformed to vegetative  
293 cells that migrate from the sediment to the upper water column. We hypothesize that the large summer bloom in 2019  
294 was the result of unusually large recruitment of akinetes in this year. (Karlsson-Elfgren et al., 2005; Karlsson-Elfgren  
295 et al., 2004). The life cycle of cyanobacteria is not a process included in the PB model (but see Hense and Beckmann  
296 (2006) and Jöhnk et al. (2011)), so increased recruitment of akinetes could explain the underestimation of the 2019  
297 summer bloom. Even the LSTM algorithms could not account for previous conditions so far back in time as to affect  
298 the formation and deposition of cyanobacteria akinetes (This may require the memory of last ice-free season). The  
299 consequent poor fit of summer bloom in 2019 partially lead to the higher *MAE* and *RMSE* in the testing dataset  
300 compared to the training dataset in all three workflows, in both GBR and LSTM models.

301 Warm winters can initiate a chain of events, i.e., shortening the ice cover duration, extending spring circulation,  
302 affected nutrients availability, and an earlier spring bloom (Adrian et al., 2006; Yang et al., 2016). According to the  
303 ice record in Lake Erken (See Fig. S1, SI), in 2020, the lake was covered by very thin ice for only 5 days, which is the  
304 shortest duration since observations were first recorded in 1954. The spring bloom in 2020 did occur earlier than other  
305 years (See Fig. S3, SI), and both ML models which considered the timing of lake ice show fairly good performance  
306 in predicting the timing and magnitude of this abnormally early spring bloom (Fig. 2, 5)

#### 307 4.1.1 Performance of Hybrid PB ML models

308 One dimensional PB hydrodynamic models can accurately simulate both water temperature profiles, and other  
309 hydrodynamic features in Lake Erken using the same forcing data that are commonly input to ML models. The hybrid  
310 model structure tested here provides a richer set of input data leading to more accurate ML predictions of algal *Chl* at  
311 little additional computational cost or data requirements. Using data from the hydrothermal PB model allowed the  
312 seasonal deepening of the thermocline, variations in the surface mixing layer depth, and upwelling events, represented  
313 by  $W_n$ , to be encoded into the ML algorithms. These factors can affect the underwater light climate, the internal loading  
314 of phosphorus and the transport of resting cyanobacteria colonies from the hypolimnion into the epilimnion favouring  
315 summer blooms of cyanobacteria (Pierson et al., 1992; Pettersson, 1998). The inclusion of these factors did increase  
316 the accuracy of the ML models, especially in the case of unusual environmental conditions (e.g. spring of 2020, Fig.  
317 2, 5) that did not frequently occur in the remaining meteorological, hydrological and biogeochemical training data.

#### 318 4.1.2 Prediction of bloom timing

319 For the purposes of water management, it may be most important to first predict the potential occurrence of a bloom,  
320 and then once underway improve predictions of its magnitude. The best model performance in predicting the timing  
321 of algal blooms, was obtained after adding hydrodynamic features derived from a PB model in workflow 3, with TPR  
322 above 45% in detecting the onset of algal bloom during 2017-2020 and a modified accuracy (Kappa) around 80 %  
323 indicated a moderate – strong level of prediction.

324 Based on our shuffling year tests of bloom timing, the GBR model showed relatively higher median TPRs than LSTM  
325 model for sample intervals less than one month. However, in some training and testing year combinations, TPRs are  
326 close to 0 % (Fig. 5), and CVs of the TPRs are highly variable, even at the original sample interval, being over 30%  
327 for GBR and over 60% for LSTM, indicating that the correct detection of algal blooms in both models are highly  
328 dependent on the years used to train the models. Thus, while the ML models can be better than the PB models at

329 predicting the onset of algal blooms, they still may not be good enough for operational forecasting. The resulting  
330 variability provided a more accurate estimate of the model performance at each down-sampled data interval and  
331 showed that increasing sample interval led to reduced performance for both ML models, in terms of *MAE*, *RMSE*, and  
332 the CV of TPR. These tests also highlighted that the performance of both ML models, especially LSTM, varied with  
333 the sampled history of events in the training period for evaluating a specific pattern of change in the testing period.  
334 We suggest that testing strategies similar to the shuffle methods used in this study are needed to accurately evaluate  
335 the expected accuracy of ML models when applied to any given site. The estimated uncertainty in shuffling training  
336 year tests (Fig. 4) and shuffling training/testing year tests (Fig. 5) can be used to better represent the uncertainty of  
337 ML derived forecasts.

#### 338 **4.2 Future applications in short-term forecasts and water management**

339 To reach the goal of incorporating ML models into operational forecasts either for short-term management support or  
340 longer-term evaluation and planning, two steps must occur. First the ML model must be developed, trained and  
341 evaluated on the water body of interest due to the unique physical characteristics and water quality dynamics in  
342 different systems. Secondly, future forcing data for the model must be obtained and integrated into a workflow that  
343 makes the future predications. In regards to the second point, a lack of frequent water monitoring (Stanley et al., 2019)  
344 is a major deterrence to applying ML models to many lakes. The data sparsity test (Fig. 5) showed that, at least for  
345 Lake Erken, the ML models can still detect the seasonal algal dynamics even for sample intervals approaching one  
346 month (Fig. S7-9). If this result holds for other lakes, the use of the two-step ML workflow could offer a method of  
347 forecasting seasonal variations in algal *Chl* even in lakes with relatively infrequent nutrient monitoring but higher  
348 frequency meteorological and hydrological data.

349 The hybrid PB/ML models have the potential to provide reasonably accurate and timely short-term algal bloom  
350 forecasts, working as part of an early-warning systems for the water resource management (Baracchini et al., 2020),  
351 and clearly have the ability to predict border seasonal variations in algal *Chl* concentration. However, since a large  
352 amount of water temperature and water quality samples are required for ML training, and since our results apply to  
353 only one well-studied lake, obtaining more datasets to test and evaluate the workflows developed here are needed.  
354 Monitoring networks (e.g., Global Lake Ecological Observatory Network [GLEON, <https://gleon.org/>]), could  
355 provide the data to allow more extensive testing and application of hybrid PB/ML models, and we are presently  
356 working in the GLEON network to test the methods developed in this paper on many other lakes.



357 **5 Code availability**

358 Model version 1.0 has been archived in Zenodo under DOI:[10.5281/zenodo.7149563](https://doi.org/10.5281/zenodo.7149563), and is available at  
359 [https://github.com/Shuqi-Lin/Erken\\_Algal\\_Bloom\\_Machine\\_Learning\\_Model.git](https://github.com/Shuqi-Lin/Erken_Algal_Bloom_Machine_Learning_Model.git).

360 **6 Data availability**

361 All data from this study have been archived with the code are also archived in Zenodo under same  
362 DOI:[10.5281/zenodo.7149563](https://doi.org/10.5281/zenodo.7149563) in the ‘training data’ folder. Here we also provide the model forcing data in the format  
363 used in the machine learning models. Data collected by the Erken laboratory, in the archived format used by the  
364 Swedish Infrastructure for Ecosystem Science (SITES) is available from the SITES data archive  
365 <https://data.fieldsites.se/portal/>

366 **7 Supplement**

367 **8 Author contribution**

368 The concept of ML model workflow was designed by SL and DP. SL developed the ML model code and performed  
369 the simulations. JM conducted the PB model simulations. SL wrote the manuscript with contributions from DP and  
370 JM.

371 **9 Competing interests**

372 The contact author has declared that neither they nor their co-authors have any competing interests.

373 **10 Acknowledgement**

374 S.L. and this study are funded by the EU and FORMAS project 2018-02771, in the frame of the collaborative  
375 international Consortium BLOOWATER (<https://www.bloowater.eu/>) financed under the ERA-NET  
376 WaterWorks2017 Cofounded Call. This ERA-NET is an integral part of the 2018 Joint Activities developed by the  
377 Water Challenges for a Changing World Joint Program Initiative (Water JPI). J.P.M. was funded by the European  
378 Union’s Horizon 2020 Research and Innovation Programme under grant agreements no. 722518 (MANTEL ITN) and  
379 101017861 (SMARTLAGOON). This study has been made possible by the Swedish Infrastructure for Ecosystem  
380 Science (SITES), in this case by data from the Erken Laboratory of Uppsala University. SITES receives funding  
381 through the Swedish Research Council under the grant no. 2017-00635.

382 **References**

- 383 Adrian, R., Wilhelm, S., and Gerten, D.: Life-history traits of lake plankton species may govern their phenological  
384 response to climate warming, *Global Change Biology*, 12, 652-661, 10.1111/j.1365-2486.2006.01125.x, 2006.
- 385 Baracchini, T., Wüest, A., and Bouffard, D.: MeteoLakes: An operational online three-dimensional forecasting  
386 platform for lake hydrodynamics, *Water Research*, 172, 115529, 10.1016/j.watres.2020.115529, 2020.
- 387 Brookes, J. D. and Carey, C. C.: Resilience to Blooms, *Science*, 334, 46-47, doi:10.1126/science.1207349, 2011.
- 388 Bruggeman, J. and Bolding, K.: A general framework for aquatic biogeochemical models, *Environmental Modelling*  
389 & Software, 61, 249-265, <https://doi.org/10.1016/j.envsoft.2014.04.002>, 2014.
- 390 Burchard, H., Bolding, K., and Villarreal, M. R.: GOTM, a General Ocean Turbulence Model: Theory,  
391 Implementation and Test Cases, European Commission. Joint Research Centre, Space Applications Institute, 103,  
392 [https://books.google.be/books/about/GOTM\\_a\\_General\\_Ocean\\_Turbulence\\_Model.html?id=zsJUHAACA AJ&redir\\_esc=y](https://books.google.be/books/about/GOTM_a_General_Ocean_Turbulence_Model.html?id=zsJUHAACA AJ&redir_esc=y), 1999.
- 394 Burford, M. A., Carey, C. C., Hamilton, D. P., Huisman, J., Paerl, H. W., Wood, S. A., and Wulff, A.: Perspective:  
395 Advancing the research agenda for improving understanding of cyanobacteria in a future of global change, *Harmful*  
396 *Algae*, 91, 101601, <https://doi.org/10.1016/j.hal.2019.04.004>, 2020.
- 397 Carey, C. C., Ibelings, B. W., Hoffmann, E. P., Hamilton, D. P., and Brookes, J. D.: Eco-physiological adaptations  
398 that favour freshwater cyanobacteria in a changing climate, *Water Research*, 46, 1394-1407,  
399 10.1016/j.watres.2011.12.016, 2012.
- 400 Elliott, J. A.: Is the future blue-green? A review of the current model predictions of how climate change could affect  
401 pelagic freshwater cyanobacteria, *Water Research*, 46, 1364-1371, 10.1016/j.watres.2011.12.018, 2012.
- 402 Friedman, J. H.: Greedy Function Approximation: A Gradient Boosting Machine, *The Annals of Statistics*, 29,  
403 1189-1232, 2001.
- 404 Hanson, P. C., Stillman, A. B., Jia, X., Karpatne, A., Dugan, H. A., Carey, C. C., Stachelek, J., Ward, N. K., Zhang,  
405 Y., Read, J. S., and Kumar, V.: Predicting lake surface water phosphorus dynamics using process-guided machine  
406 learning, *Ecological Modelling*, 430, 109136, 10.1016/j.ecolmodel.2020.109136, 2020.
- 407 Hense, I. and Beckmann, A.: Towards a model of cyanobacteria life cycle—effects of growing and resting stages on  
408 bloom formation of N<sub>2</sub>-fixing species, *Ecological Modelling*, 195, 205-218,  
409 <https://doi.org/10.1016/j.ecolmodel.2005.11.018>, 2006.
- 410 Hochreiter, S. and Schmidhuber, J.: Long Short-Term Memory, *Neural Computation*, 9, 1735-1780,  
411 10.1162/neco.1997.9.8.1735, 1997.
- 412 Huisman, J., Codd, G. A., Paerl, H. W., Ibelings, B. W., Verspagen, J. M. H., and Visser, P. M.: Cyanobacterial  
413 blooms, *Nature Reviews Microbiology*, 16, 471-483, 10.1038/s41579-018-0040-1, 2018.
- 414 Jia, X., Willard, J., Karpatne, A., Read, J., Zwart, J., Steinbach, M., and Kumar, V.: Physics Guided RNNs for  
415 Modeling Dynamical Systems: A Case Study in Simulating Lake Temperature Profiles, in: *Proceedings of the 2019*  
416 *SIAM International Conference on Data Mining (SDM)*, 558-566, 2019.
- 417 Jimeno-Sáez, P., Senent-Aparicio, J., Cecilia, J. M., and Pérez-Sánchez, J.: Using Machine-Learning Algorithms for  
418 Eutrophication Modeling: Case Study of Mar Menor Lagoon (Spain), *International Journal of Environmental*  
419 *Research and Public Health*, 17, 1189, 2020.
- 420 Jöhnk, K. D., Brüggemann, R., Rucker, J., Luther, B., Simon, U., Nixdorf, B., and Wiedner, C.: Modelling life cycle  
421 and population dynamics of Nostocales (cyanobacteria), *Environmental Modelling & Software*, 26, 669-677,  
422 <https://doi.org/10.1016/j.envsoft.2010.11.001>, 2011.
- 423 Karlsson-Elfgren, I., Hyenstrand, P., and Riydin, E.: Pelagic growth and colony division of *Gloeotrichia echinulata*  
424 in Lake Erken, *Journal of Plankton Research*, 27, 145-151, DOI 10.1093/plankt/fbh165, 2005.
- 425 Karlsson-Elfgren, I., Rengefors, K., and Gustafsson, S.: Factors regulating recruitment from the sediment to the  
426 water column in the bloom-forming cyanobacterium *Gloeotrichia echinulata*, *Freshwater Biology*, 49, 265-273, DOI  
427 10.1111/j.1365-2427.2004.01182.x, 2004.
- 428 Marcé, R., George, G., Buscarinu, P., Deidda, M., Dunalska, J., de Eyto, E., Flaim, G., Grossart, H.-P., Istvanovics,  
429 V., Lenhardt, M., Moreno-Ostos, E., Obrador, B., Ostrovsky, I., Pierson, D. C., Potužák, J., Poikane, S., Rinke, K.,  
430 Rodríguez-Mozaz, S., Staehr, P. A., Šumberová, K., Waajen, G., Weyhenmeyer, G. A., Weathers, K. C., Zion, M.,  
431 Ibelings, B. W., and Jennings, E.: Automatic High Frequency Monitoring for Improved Lake and Reservoir  
432 Management, *Environmental Science & Technology*, 50, 10780-10794, 10.1021/acs.est.6b01604, 2016.
- 433 McHugh, M. L.: Interrater reliability: the kappa statistic, *Biochemia medica*, 22, 276-282, 2012.
- 434 Mellios, N., Moe, S. J., and Laspidou, C.: Machine Learning Approaches for Predicting Health Risk of Cyanobacterial  
435 Blooms in Northern European Lakes, *Water*, 12, 1191, 2020.

437 Mesman, J. P., Ayala, A. I., Goyette, S., Kasparian, J., Marcé, R., Markensten, H., Stelzer, J. A. A., Thayne, M. W.,  
438 Thomas, M. K., Pierson, D. C., and Ibelings, B. W.: Drivers of phytoplankton responses to summer wind events in a  
439 stratified lake: A modeling study, *Limnology and Oceanography*, 67, 856-873, <https://doi.org/10.1002/lno.12040>,  
440 2022.

441 Moras, S., Ayala, A. I., and Pierson, D. C.: Historical modelling of changes in Lake Erken thermal conditions,  
442 *Hydrology and Earth System Sciences*, 23, 5001-5016, 2019.

443 Nelson, N. G., Muñoz-Carpena, R., Philips, E. J., Kaplan, D., Sucsy, P., and Hendrickson, J.: Revealing Biotic and  
444 Abiotic Controls of Harmful Algal Blooms in a Shallow Subtropical Lake through Statistical Machine Learning,  
445 *Environmental Science & Technology*, 52, 3527-3535, 10.1021/acs.est.7b05884, 2018.

446 Paerl, H. W.: Nuisance phytoplankton blooms in coastal, estuarine, and inland waters<sup>1</sup>, *Limnology and*  
447 *Oceanography*, 33, 823-843, 10.4319/lno.1988.33.4part2.0823, 1988.

448 Paerl, H. W. and Huisman, J.: Blooms Like It Hot, *Science*, 320, 57-58, doi:10.1126/science.1155398, 2008.

449 Persson, I. and Jones, I. D.: The effect of water colour on lake hydrodynamics: a modelling study, *Freshwater*  
450 *Biology*, 53, 2345-2355, <https://doi.org/10.1111/j.1365-2427.2008.02049.x>, 2008.

451 Pettersson, K.: The Availability of Phosphorus and the Species Composition of the Spring Phytoplankton in Lake  
452 Erken, *Internationale Revue der gesamten Hydrobiologie und Hydrographie*, 70, 527-546,  
453 10.1002/iroh.19850700407, 1985.

454 Pettersson, K.: Mechanisms for internal loading of phosphorus in lakes, *Hydrobiologia*, 373, 21-25,  
455 10.1023/A:1017011420035, 1998.

456 Pettersson, K., Grust, K., Weyhenmeyer, G., and Blenckner, T.: Seasonality of chlorophyll and nutrients in Lake  
457 Erken – effects of weather conditions, *Hydrobiologia*, 506, 75-81, 10.1023/B:HYDR.0000008582.61851.76, 2003.

458 Peretyatko, A., Teissier, S., De Backer, S., and Triest, L.: Classification trees as a tool for predicting cyanobacterial  
459 blooms, *Hydrobiologia*, 689, 131-146, 10.1007/s10750-011-0803-4, 2012.

460 Pierson, D. C., Pettersson, K., and Istvanovics, V.: Temporal changes in biomass specific photosynthesis during the  
461 summer: regulation by environmental factors and the importance of phytoplankton succession, *Hydrobiologia*, 243,  
462 119-135, 10.1007/BF00007027, 1992.

463 Read, J. S., Hamilton, D. P., Jones, I. D., Muraoka, K., Winslow, L. A., Kroiss, R., Wu, C. H., and Gaiser, E.:  
464 Derivation of lake mixing and stratification indices from high-resolution lake buoy data, *Environmental Modelling*  
465 *& Software*, 26, 1325-1336, 10.1016/j.envsoft.2011.05.006, 2011.

466 Read, J. S., Jia, X., Willard, J., Appling, A. P., Zwart, J. A., Oliver, S. K., Karpatne, A., Hansen, G. J. A., Hanson, P.  
467 C., Watkins, W., Steinbach, M., and Kumar, V.: Process-Guided Deep Learning Predictions of Lake Water  
468 Temperature, *Water Resources Research*, 55, 9173-9190, 10.1029/2019WR024922, 2019.

469 Rousso, B. Z., Bertone, E., Stewart, R., and Hamilton, D. P.: A systematic literature review of forecasting and  
470 predictive models for cyanobacteria blooms in freshwater lakes, *Water Research*, 182, 115959,  
471 10.1016/j.watres.2020.115959, 2020.

472 Recknagel, F., Fukushima, T., Hanazato, T., Takamura, N., and Wilson, H.: Modelling and prediction of phyto- and  
473 zooplankton dynamics in Lake Kasumigaura by artificial neural networks, *Lakes & Reservoirs: Science, Policy and*  
474 *Management for Sustainable Use*, 3, 123-133, 10.1111/j.1440-1770.1998.tb00039.x, 1998.

475 Reichwaldt, E. S. and Ghadouani, A.: Effects of rainfall patterns on toxic cyanobacterial blooms in a changing  
476 climate: Between simplistic scenarios and complex dynamics, *Water Research*, 46, 1372-1393,  
477 10.1016/j.watres.2011.11.052, 2012.

478 Richardson, J., Miller, C., Maberly, S. C., Taylor, P., Globevnik, L., Hunter, P., Jeppesen, E., Mischke, U., Moe, S.  
479 J., Pasztaleniec, A., Søndergaard, M., and Carvalho, L.: Effects of multiple stressors on cyanobacteria abundance  
480 vary with lake type, *Global Change Biology*, 24, 5044-5055, 10.1111/gcb.14396, 2018.

481 Rousso, B. Z., Bertone, E., Stewart, R., and Hamilton, D. P.: A systematic literature review of forecasting and  
482 predictive models for cyanobacteria blooms in freshwater lakes, *Water Research*, 182, 115959,  
483 10.1016/j.watres.2020.115959, 2020.

484 Stanley, F. K. T., Irvine, J. L., Jacques, W. R., Salgia, S. R., Innes, D. G., Winkvist, B. D., Torr, D., Brenner, D. R.,  
485 and Goodarzi, A. A.: Radon exposure is rising steadily within the modern North American residential environment,  
486 and is increasingly uniform across seasons, *Scientific Reports*, 9, 18472, 10.1038/s41598-019-54891-8, 2019.

487 Watson, S. B., Miller, C., Arhonditsis, G., Boyer, G. L., Carmichael, W., Charlton, M. N., Confesor, R., Depew, D.  
488 C., Höök, T. O., Ludsin, S. A., Matisoff, G., McElmurry, S. P., Murray, M. W., Peter Richards, R., Rao, Y. R.,  
489 Steffen, M. M., and Wilhelm, S. W.: The re-eutrophication of Lake Erie: Harmful algal blooms and hypoxia,  
490 *Harmful Algae*, 56, 44-66, <https://doi.org/10.1016/j.hal.2016.04.010>, 2016.

491 Wei, B., Sugiura, N., and Maekawa, T.: Use of artificial neural network in the prediction of algal blooms, *Water*  
492 *Research*, 35, 2022-2028, 10.1016/S0043-1354(00)00464-4, 2001.

493 Wilson, H. L., Ayala, A. I., Jones, I. D., Rolston, A., Pierson, D., de Eyto, E., Grossart, H.-P., Perga, M.-E.,  
494 Woolway, R. I., and Jennings, E.: Variability in epilimnion depth estimations in lakes, *Hydrology and Earth System*  
495 *Sciences*, 24, 5559-5577, 10.5194/hess-24-5559-2020, 2020.  
496 Xiao, X., He, J., Huang, H., Miller, T. R., Christakos, G., Reichwaldt, E. S., Ghadouani, A., Lin, S., Xu, X., and Shi,  
497 J.: A novel single-parameter approach for forecasting algal blooms, *Water Research*, 108, 222-231,  
498 10.1016/j.watres.2016.10.076, 2017.  
499 Yang, Y., Stenger-Kovács, C., Padisák, J., and Pettersson, K.: Effects of winter severity on spring phytoplankton  
500 development in a temperate lake (Lake Erken, Sweden), *Hydrobiologia*, 780, 47-57, 10.1007/s10750-016-2777-8,  
501 2016.