# Prediction of algal blooms via data-driven machine learning models:

# An evaluation using data from a well monitored mesotrophic lake

Shuqi Lin[1*], Donald C. Pierson[1], Jorrit P. Mesman [1,2]

[1]Erken Laboratory and Limnology Department, Uppsala University, Uppsala, Sweden

[2]Département F.-A. Forel des sciences de l'environnement et de l'eau, Université de Genève, Genève,

Switzerland

*Correspondence to*: Shuqi Lin (Shuqi.lin@ebc.uu.se)

**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 blooms, methods to forecast their timing and magnitude are needed. However, the factors regulating algal blooms are complex, variable and site-specific, often involving high-order interactions of environmental factors and 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 (Elliott, 2012).

32  With the proliferation of lake monitoring data (Marcé et al., 2016), data-driven machine learning (ML) approaches

33  have been applied, as an alternative to PB models for bloom prediction (Rousso et al., 2020). Previously applied

34  ML models, including Random Forest (Recknagel et al., 1998), Support Vector Machine (Jimeno-Sáez et al.,

35  2020), and Artificial Neural Network (Xiao et al., 2017; Nelson et al., 2018; Wei et al., 2001), can improve

36  predictions of the timing and seasonality of algal *Chl* pattern, apparently by accounting for complexity that is

37  difficult to encode within the framework of a PB model.

38  In this study, we propose a two-step ML approach for predicting algal dynamics that: first estimates lake nutrient

39  concentrations which often have limited observations and secondly predicts variations in algal *Chl* using these

40  pre-generated nutrient concentrations combined with other observed environmental factors that are collected at

41  higher frequency. We also test a simple hybrid model architecture that by adding hydrodynamic features derived

42  from the PB model into the training features of the two-step ML approach, allowing us to include additional

43  information describing physical lake processes expected to affect variations in algal growth and succession in the

44  machine learning prediction.

45  We applied the above workflows to predict changing *Chl* concentration, as a proxy for the occurrence of algal

46  blooms, via Gradient Boost Regressor (GBR) and Long Short-term Memory network (LSTM). Two shuffling

47  year tests were conducted. One assessed the uncertainty of ML models in predicting *Chl* during the same two-

48  year period and the other evaluated the sensitivity of ML accuracy to various training/testing year combinations

49  and lake nutrient sampling intervals. Model performance and potential applications in algal bloom forecasting are
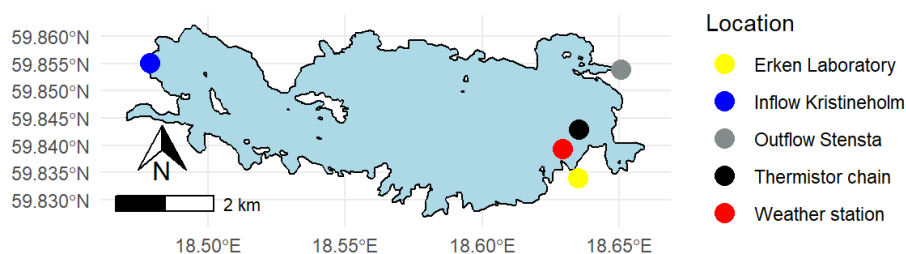
50  discussed.

51  **2 Methods**

52  **2.1 Study site**

53  The study site, Lake Erken, is a mesotrophic lake located in east-central Sweden, that has a surface area of 24

54  km$^2$, a maximum depth of 21 m and an average retention time of 7 years. The lake is dimictic with seasonal

55  stratification commonly beginning in May-June and ending in August-September. The onset of ice cover usually

56  beings in December-February and the loss of ice occurs in Mar-April (Persson and Jones, 2008). Located near the

57  Baltic coast, Lake Erken is wind exposed, and susceptible to periodic wind-induced turbulent mixing.

58　Changes in algal *Chl* in Lake Erken have a typical seasonal pattern, with spring and summer peaks in concentration

59　(Pettersson et al., 2003). Spring blooms are dominated by dinoflagellates and diatoms (Pettersson, 1985), and

60　initiated by overwinter species from the last autumn (Yang et al., 2016). Cyanobacteria dominate summer peaks

61　in *Chl*, given that they can optimize their vertical position in regarding to nutrients and light (Paerl, 1988; Pierson

62　et al., 1992).



63

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

65　**2.2 Data**

66　Lake Erken has a long running automated monitoring program that provides hourly meteorological data, water

67　temperature profiles between 0.5 and 15 m at 0.5 m intervals and the flow from the inflow and outflow (Fig.1). A

68　manual sampling program collects samples during ice-free time at 5-7 days intervals for all major nutrient

69　concentrations (e.g., $NO_X$, $NH_4$, $PO_4$, Total P, Si, etc.), dissolved oxygen ($O_2$), and *Chl* concentration. The timing

70　of the onset and loss of ice cover are also monitored yearly by the lab. More detailed information on the sampling

71　program is in Supporting Information (See Text S1) and Moras et al. (2019).

72　**2.3 Modelling Methods**

73　2.3.1 Process-based (PB) lake model

74　In this study, a PB hydrodynamic lake model, GOTM (General Ocean Turbulence Model) (Burchard et al., 1999),

75　was used to generate water temperature profiles, and other hydrodynamic metrics. GOTM also served as the

76　foundation of water quality simulations made with the SELMAPROTBAS model (Mesman et al., 2022) that is

77　coupled to GOTM through the Framework for Aquatic Biogeochemical Models FABM (Bruggeman and Bolding,

78　2014).

79　2.3.2 Data-driven machine learning (ML) models

80　Two ML models were evaluated in this study. Gradient Boosting Regressor (GBR) which iteratively generates an

81　ensemble of estimator trees with each tree improving upon the performance of the previous (Friedman, 2001), and

82  Long short-term memory (LSTM) networks which is built for sequential and timeseries modelling (Hochreiter

83  and Schmidhuber (1997), See Fig. S2, SI). The hyperparameter settings in both ML models can be found in

84  Supporting Information (See Text S2). Both ML models are built on Python using the Scikit-Learn (https://scikit-

85  learn.org/stable/, last access: September, 2021) and TensorFlow (https://www.tensorflow.org/, last access:

86  September, 2021) libraries.

87  **2.4 Design of predictive workflows and shuffling year data sparsity tests**

88  In this study, we tested three workflows using a dataset split for training (years 2004-2016) and testing (years

89  2017-2020). In all three workflows, a 5-fold cross-validation using the training dataset was used to optimize the

90  hyperparameters in the ML models. Workflow 1 directly predicts *Chl* concentration based on available

91  environmental observations (See SI, Table S1). The training and testing datasets were limited by the frequency of

92  lake nutrient observations which resulted in 5-7 days gap between data points. The time step of LSTM was set to

93  1, that is, the environmental factors on the target date and previous observation date, which may be 5-7 days ago,

94  were used to train the model and make predictions.

95  In workflow 2 and 3, a two-step approach was applied (Table S1). Daily measurements of physical factors were

96  used to pre-generate daily variations in lake nutrients via separate ML models, and the ML models were trained

97  at a daily time step using the measured environmental factors and pre-generated nutrient concentrations. The time

98  step of LSTM was then set to 7 days.

99  In workflow 3, three hydrodynamic features, i.e., mixing layer depth ($z_e$), Wedderburn number ($W_n$), and the

100  seasonal thermocline depth (*thermD*), derived from the GOTM model were regarded as daily training features in

101  the two-step ML approach. The definitions and calculations of these features are explained in SI (2.5 Feature

102  selection and processing for ML models, Text S3)

103  Following the two-step approach and using workflow 3, we set up two tests. (1) To assess the uncertainty induced

104  by variations in the data used to train the ML models, we shuffled the training years, randomly taking 13 years

105  out of 2004-2018 dataset 30 times, and tested the model predictions of *Chl* during 2019-2020. And, (2) to test if

106  the workflow could be used for other water systems which may have less frequent lake nutrient monitoring data,

107  we conducted a data sparsity test that evaluated the sensitivity of models to the lake nutrient and *Chl* sampling

108  interval. For this test the lake nutrient and *Chl* concentration observations in training dataset was down-sampled

109  to a 7-day, 14-day, 21- day, 28-day, and 35-day sampling interval. Then for each sampling interval using the 2004-

110  2020 dataset, *Chl* was predicted for different consecutive 4-year periods when the ML models were trained by the

111    remaining 13 years of data. Data shuffling was conducted 13 times so that every 4-year period in our dataset was

112    tested.

113    **2.5 Feature selection and processing for ML models**

114    The feature selection process is based on some a priori knowledge of the underlying phenomena related to algal

115    blooms. All workflows made use of the daily automated monitoring data. In addition, the temperature difference

116    ($\Delta T$) between surface water (averaged over the upper 3 m) and bottom water (15 m) was also used to represent

117    the thermal structure of the lake., and the duration of ice cover in the previous winter, and the number of days

118    from ice-off date were used.

119    In workflow 2 and 3 nutrients are predicted sequentially, with each pre-generated nutrient predictions included in

120    the training data of the next nutrient prediction (Table S1). Workflow 3 added $z_e$, computed using the GOTM

121    simulated vertical eddy diffusivity ($K_z$) profiles, *thermD*, estimated using Lake Analyzer (Read et al., 2011) based

122    on GOTM simulated temperature profile, and $W_n$, a dimensionless parameter measuring the balance between wind

123    stress and the pressure gradient resulting from the slope of the interface (See Text S3, SI), as additional daily

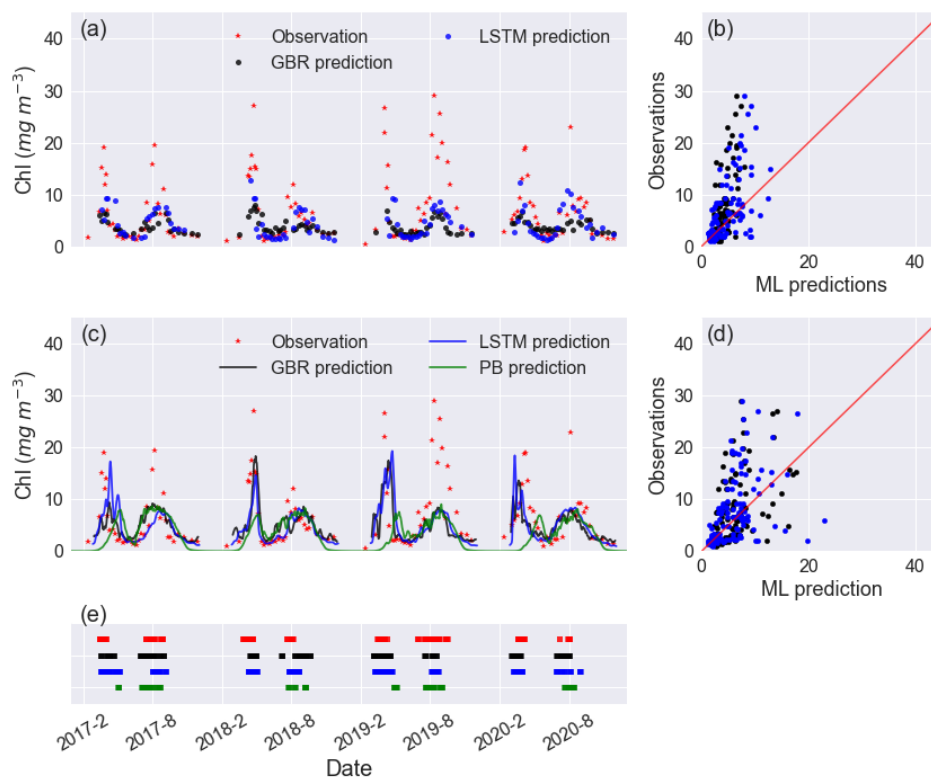124    training features.

125    **2.6 Evaluating metrics**

126    Model performance was evaluated by comparing the simulated and measured *Chl* concentrations, and by

127    calculating the mean absolute error (*MAE*), root means square error (*RMSE*), and correlation coefficient ($R^2$). To

128    evaluate the accuracy of the model in detecting the onset of an algal bloom, we calculated a confusion matrix in

129    workflows 2 and 3, where the observations were linearly interpolated to daily values, and predicted daily *Chl*

130    concentration were smoothed with a 7-day rolling mean. Using these data, the onset of a bloom was categorized

131    as occurring when the daily change of *Chl* (*ΔChl*) exceed a threshold, 0.35 mg m$^{-3}$ day$^{-1}$. This works well in Lake

132    Erken where *Chl* concentrations are frequently monitored (near weekly), and the linear interpolation can be

133    expected to be reasonably representative of the *Chl* concentrations between measured samples. Considering the

134    randomization in the ML models, we also add a 3-day window on the bloom onset prediction, that is, we

135    considered the prediction of a bloom valid if the measured data suggested a bloom the day before or after the

136    simulated onset. We used the True Positive Rate (TPR), False Positive Rate (FPR), and modified accuracy Kappa

137    (Mchugh, 2012) to identify the potential of ML models to correctly capture the algal bloom onset (See Table S2,

138    SI). A model with 100% TPR, 0% FPR, and 100% Kappa would constitute a perfect fit.

139    **3 Results**

140    **3.1 Workflow 1: Direct prediction based on observations**

141    In workflow 1, both GBR and LSTM clearly reproduced spring and summer blooms (Fig. 2a) but underestimated

142    the intensity of blooms (Fig. 2a, b). Neither ML model captured the extraordinarily high *Chl* (~15-30 mg m$^{-3}$) in

143    the summer of 2019. The cross-validation on training dataset (See Table S3, SI) shows what appears to be

144    overfitting issue in both models, with somewhat higher *RMSE* and *MAE* in the testing dataset than the mean values

145    in the training dataset. The achieved accuracy of models is attributed to the daily availability of physical inputs,

146    and the fact that in Lake Erken water samples are collected frequently at 5-7 days intervals. Workflow 1 may be

147    most valuable in reconstructing previous variations in algal *Chl*, filling the gaps between measured *Chl*

148    observations and feature importance ranking (See Fig. S4, SI). But when using this workflow future forecasts will

149    be limited by the absence of future nutrient data.



150

151    **Figure 2.** Timeseries of observed and predicted *Chl* from GBR and LSTM models in (a) workflow 1 and (c)

152    workflow 3, and the corresponding scatter plots of observations vs ML predictions of *Chl* in workflow 1 and

153    workflow 3 are shown in panels (b) and (d), with the black and blue dots/lines representing the predictions from

154 GBR and LSTM, respectively. Panel (e) shows the observed and predicted algal bloom onsets in 2017-2020 using

155 the same color coding as the previous panels. Results from the PB model simulation in Mesman et al. (2022) are

156 also shown in (c) and (e).

157 **3.2 Workflow 2: Two-step ML models based on pre-generated daily nutrients and observed physical**

158 **factors**

159 As in workflow 1, both ML models in workflow 2 suffered from overfitting with higher *MAE*, *RMSE*, and lower

160 $R^2$ in testing datasets than training datasets (See SI, Table S3).

161 Overall, both GBR and LSTM showed slightly higher *MAE* (4.22 mg m$^{-3}$ vs. 3.87 mg m$^{-3}$) and *RMSE* (6.27 mg

162 m$^{-3}$ vs. 6.00 mg m$^{-3}$) when compared to workflow 1 (Table 1). But they also showed improved performance in

163 terms of capturing the peak values of *Chl* during spring blooms (Fig. 2, Fig. S5, SI). Both workflows outperformed

164 the SELMAPROTBAS PB model in simulating concentrations of lake nutrients (See Fig. S6, SI). The ML models

165 were more accurate in predicting the low values of NO$_X$ and peak values of PO$_4$ and Total P. However, both ML

166 models and the PB model failed in predicting the extremely high values of measured lake nutrients, such as the

167 autumn peak of NH$_4$ in 2017 (Fig. S6e) and the spring peak of O$_2$ in 2018 (Fig. S6c), Thus, higher workflow 2

168 *MAE* and *RMSE* (Table 1) are presumably due to the inaccuracies in the pre-generated nutrient training data, but

169 the improved daily predictions that better capture the bloom events, overshadow these flaws.

170 **Table 1.** Comparisons of model performance during the testing period based on *RMSE*, *MAE*, and *R2*. The unit

171 of *Chl* is mg m$^{-3}$.

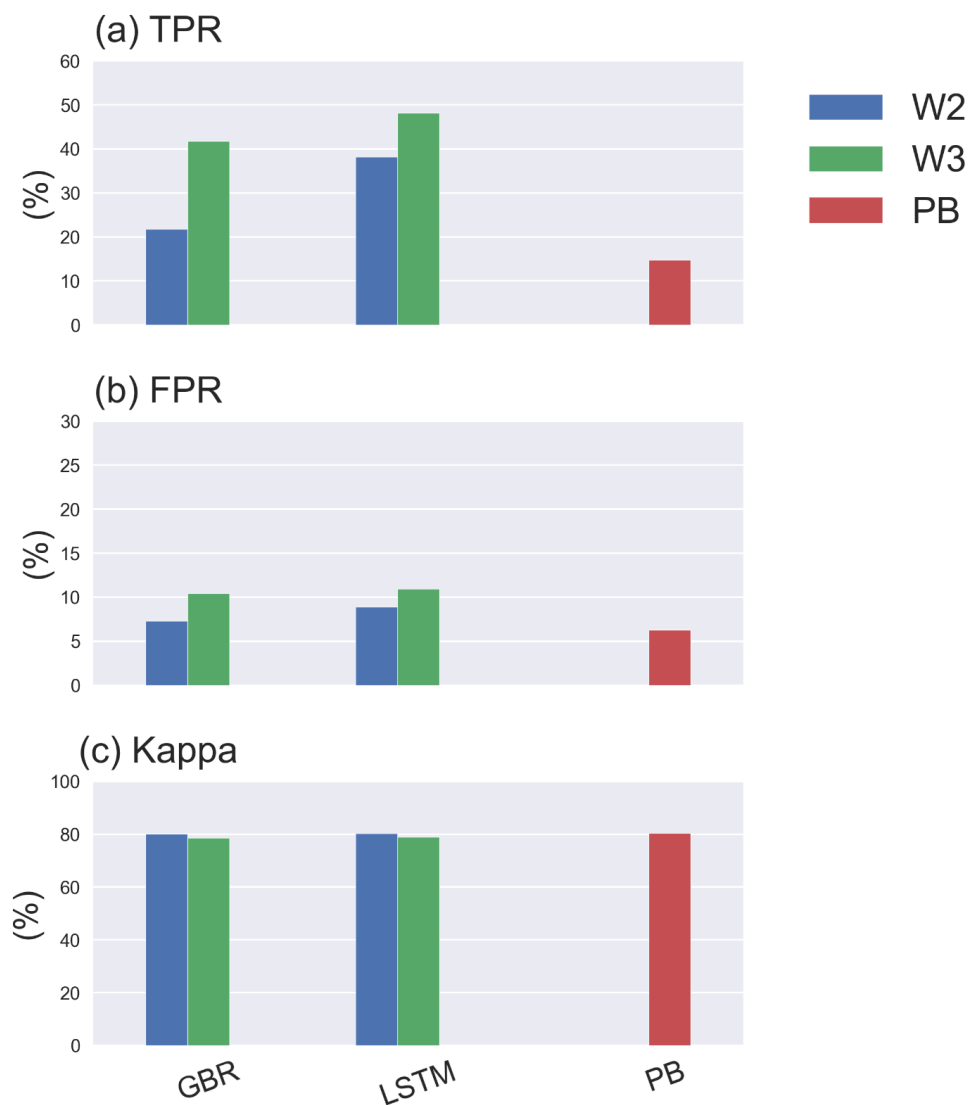| Model | PB | ML-workflow 1 | | ML-workflow 2 | | ML-workflow 3 | |
|---|---|---|---|---|---|---|---|
| | | GBR | LSTM | GBR | LSTM | GBR | LSTM |
| *RMSE* | 7.18 | 5.77 | **5.64** | 6.27 | 6.00 | 5.94 | 5.81 |
| *MAE* | 4.77 | **3.55** | 3.58 | 4.22 | 3.87 | 3.99 | 3.71 |
| *R2* | -0.25 | 0.13 | **0.20**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

172

173 **3.3 Workflow 3: based on workflow 2, and including hydrodynamic training features derived from the**

174 **GOTM model.**

175 Including hydrodynamic training information in workflow 3 did not significantly improve in lake nutrient

176 predictions compared to workflow 2 (See Fig. S6), and when using workflow 3 both ML models showed

177 comparable performance in *Chl* predictions compared to workflow 1. However, the predictions of the spring

178 bloom in all years improved compared to workflows 1 and 2, in terms of the magnitude and timing of spring

179 bloom (Fig. 2e). This was the case in 2019-2020 (Fig. 2a) which was an abnormally warm winter with only 5 days

180   ice cover, and had an unusually early spring algal bloom. Both workflow 2 and 3 did not capture the extremely

181   intensive bloom (with peak values closed to 30 mg m$^{-3}$) in summer of 2019, and neither did the PB model.

182   Furthermore, adding hydrodynamic features derived from PB model improved predictions of the onset of algal

183   blooms (Fig. 2e and 4), with the overall TPR increasing by 15 % and 5 %, FPR increasing around 5% and 3 % in

184   GBR and LSTM models, respectively. Compared with the PB model which showed lower TPR (15%) and FPR

185   (6%), ML models are more likely to predict algal bloom at the correct time. However, the concomitant higher

186   FPRs indicating an incorrect warning of algal bloom is also more likely to occur in the ML models, since the PB

187   model is more like to miss the bloom entirely. The Kappa values of both ML models and the PB model are close

188   to 80%, showing that all models simulated the entire period (blooms and the periods between blooms) to a
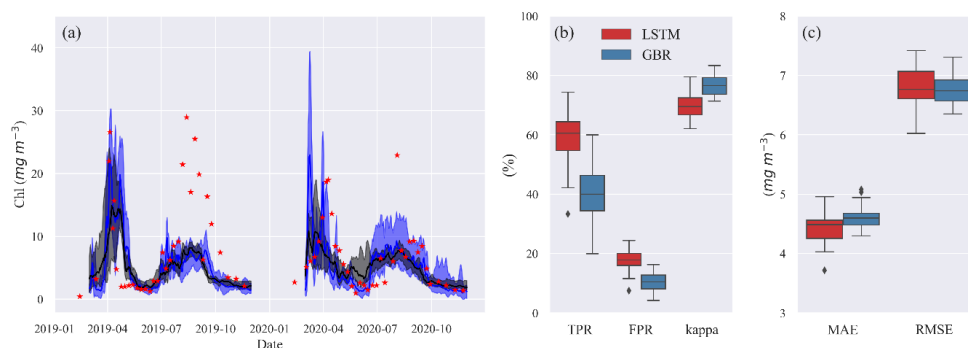
189   moderate-strong level (Mchugh, 2012).

Geoscientific
Model Development
Discussions



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

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

The results presented so far are based on a typical strategy of training ML models for a historical period in this case 2004-2016 and then accessing model performance in a second period between 2017-2020. The accuracy of the model predictions will to some extent be related to the range and variability in the training data. To evaluate the importance of this we randomly removed two years from a 2004-2018 training dataset, and made 30 different predictions of *Chl* during 2019-2020 when the models had difficulties predicting spring and summer blooms (Fig 5). When trained with the various shuffled combinations, both ML models were capable of reproducing the

199    seasonal variations in algal *Chl* with a 4.5 % and 5.8 % coefficient of variation (CV) in *MAE*, and a 24.0 % and

200    16.4 % CV in TPR of GBR and LSTM, respectively (See Table S4, SI). This provides an indication of the

201    uncertainty that may arise as a consequence of differences in the training datasets used for in our workflows. And,

202    it also shows that even a relatively long training period of 13 years can not totally capture the system behaviour

203    in such a way as to lead to nearly similar bloom predictions.

204    Although none of the model runs captured the intensive summer bloom in 2019, the spring bloom in both years

205    was well represented, especially by LSTM, in terms of timing and magnitude.



206

207    **Figure 4.** (a) Timeseries of observed (red stars) and predicted *Chl* from GBR (black) and LSTM (blue) models in

208    shuffling training year test. The shades represent the range between minimum and maximum prediction, and the

209    solid lines represent the median prediction.

210

211    Despite comparable *RMSE* and *MAE* in LSTM and GBR (Fig. 4c), both higher TPRs (with median of 60%) and

212    FRPs (with median of 18%) in LSTM indicate that the LSTM was more aggressive in making algal bloom

213    predictions. The GBR model's apparent advantage in FPRs (with median 10%) is largely the result of it making

214    a lower number of bloom predictions since the low concentrations between spring and summer blooms in 2020

215    was not well represented (Fig. 4b).

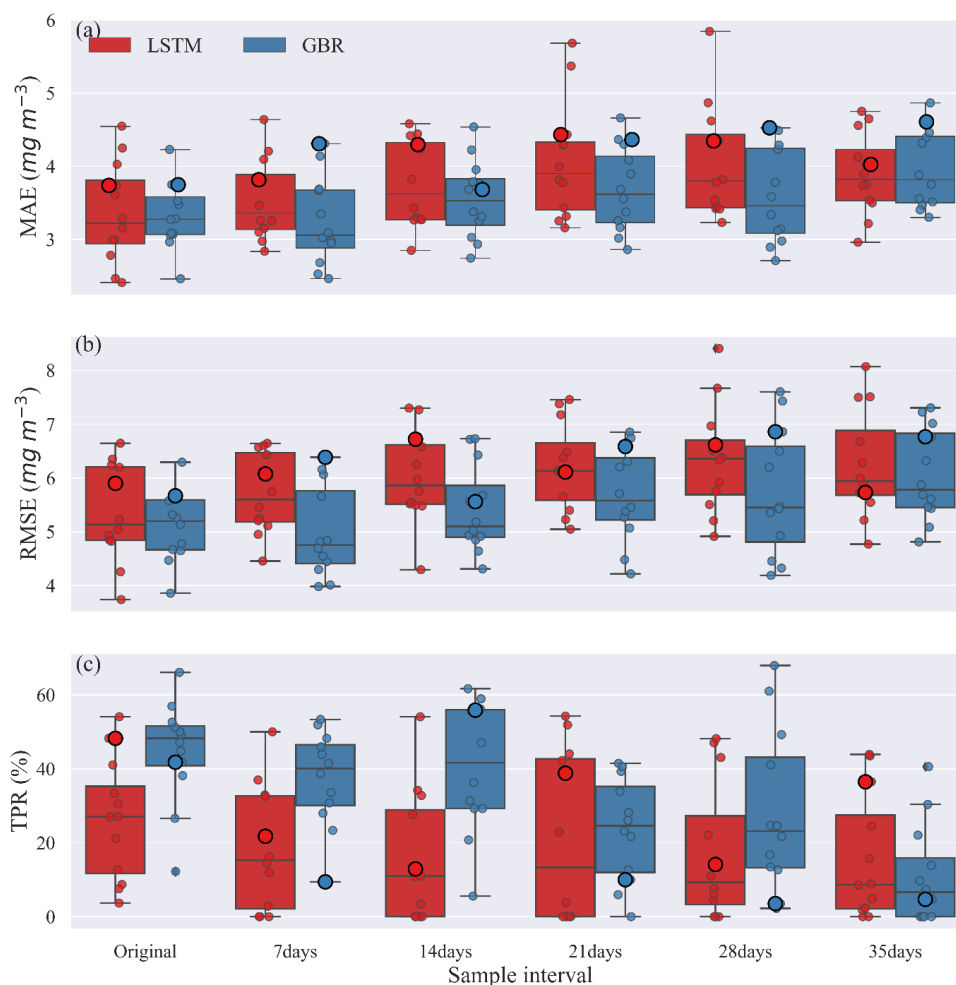216    **3.5 Shuffling years data sparsity test**

217    To examined the possible use of workflow 3 when data are less frequently available, lake nutrient and *Chl* data

218    were down-sampled so that the effects of sampling frequency on model predictions could be evaluated. Each

219    down-sampled dataset was also rearranged into 13 different 13-year training periods and 4-year testing periods.

220    The variability in predictions provided a measure of model performance and uncertainty. Fig. 5 shows the

221    uncertainty in model predictions as a consequence of the chosen sampling intervals.

222    The *MAE*s and *RMSE*s of both GBR and LSTM models tended to increase with the longer sample intervals. The

223    median *MAE* was always slightly higher for the LSTM model except when trained with original dataset (Fig. 5a).

224 While our initial evaluation of TPR using 2017-2020 as the testing period and 2004-2016 as the training period

225 suggested the LSTM model was more accurate in turns of detection of algal bloom onsets (Fig. 3), Fig. 5c showed

226 the median TPR of GBR model over was over 50%, higher than LSTM model. This can be explained by the fact

227 that the 2017-2020 period as in Fig. 3 and shown as large points in Fig. 5 was unusually difficult for GBR to

228 simulate. Consequently, even though the GBR model usually performs better in Fig. 5c the testing period chosen

229 for use in Fig. 3, showed the opposite result. This illustrates the importance of the sequence of training and testing

230 years for evaluating model performance.

231 For the first three sampling intervals the GBR model clearly had better TPR values than the LSTM model. The

232 median TPRs of GBR model started to drop below 30% once the sample interval reached 21 days. For LSTM,

233 medium TPRs remained lower than 30%, for all sampling intervals but also showed a much wider range of

234 variability (Table S5) dependent on the training and tested datasets used. In general, both models preformed best

235 at the original and 7-day sampling interval, but then showed slightly worse performance that was consistent up to

236 a sample interval of 21 days. In terms of the errors evaluated over the entire 4-year testing period (Fig. 5a, b) the

237 GBR model had lower errors and therefore, better, predicted the seasonal variations of *Chl* concentration. The

238 timeseries comparison of observed and predicted *Chl* from this shuffling year data sparsity test can be found in SI

239 (Fig. S7-9).

**Figure 5.** Comparisons of (a) *MAE*, (b) *RMSE*, and (c) TPR between GBR and LSTM under various sample intervals. Circles along the box show the result from every training and testing years combination and the bigger circles represent 2004-2016 training and 2017-2020 testing years combination as was used in Fig. 2.

**4 Discussion**

**4.1 Performance of ML models**

In three workflows, the ML models successfully reproduced the *Chl* seasonal patterns, capturing the spring and summer bloom events, with lower averaged *RMSE*s and *MAE*s than PB model simulations that was previously calibrated for use in Lake Erken. Workflow 1 which predicted *Chl* based on all available environmental factors including lake nutrient observations showed that both ML models can reproduce the seasonal dynamics of algal *Chl* with promising accuracy (*MAE* = 3.55 and 3.58 mg m$^{-3}$, *RMSE* = 5.77 and 5.64 mg m$^{-3}$ and $R^2$ = 0.13 and

251    0.20) via the direct input of available environmental observations. These ML models can be applied to reconstruct

252    past patterns of algal *Chl*, fill the gaps between measured *Chl* observations, and interpret the mechanisms that

253    drive phytoplankton dynamics. Workflows 2 and 3 adopted a two-step approach, first using separate ML models

254    to estimating daily changes in lake nutrient concentration, and in Workflow 3 also including PB model derived

255    physical factors as training features of the algal ML model. These two workflows allowed daily predictions of

256    changes in algal *Chl* concentration using both observations and pre-generated lake nutrient concentrations at a

257    consistent daily time step, and achieved comparable accuracy in *Chl* prediction to workflow 1, demonstrating the

258    potential for making daily forecasts without measured nutrient observations.

259    However, there was overfitting issues in all three workflows, in both GBR and LSTM models, indicated by higher

260    *MAE* and *RMSE* in the testing dataset compared to the training dataset especially for GBR (Table S3).

261    The one clear failure of both the ML and PB based model predictions was during July-August 2019, *Chl*

262    concentrations in integrated samples collected between the surface and 6-12 m exceeded 20 mg m$^{-3}$ over a 5-week

263    period. Neither the PB model nor ML models captured this unusually persistent bloom (Fig. 2, Fig. S3, SI). At

264    this time the phytoplankton were dominated by the cyanobacteria Gloeotrichia and Anabaena, that form a resting

265    akinete life stage at the end of their yearly bloom, which can initiate the following year's bloom as they are

266    transformed to vegetative cells that migrate from the sediment to the upper water column. We hypothesize that

267    the large summer bloom in 2019 was the result of unusually large recruitment of akinetes in this year. (Karlsson-

268    Elfgren et al., 2005; Karlsson-Elfgren et al., 2004). The life cycle of cyanobacteria is not a process included in the

269    PB model (but see Hense and Beckmann (2006) and Jöhnk et al. (2011)), so increased recruitment of akinetes

270    could explain the underestimation of the 2019 summer bloom. Even the LSTM algorithms could not account for

271    previous condition so far back in time as to affect the formation and deposition of cyanobacteria akinetes.

272    Warm winters can initiate a chain of events, i.e., shortening the ice cover duration, extending spring circulation,

273    affected nutrients availability, and an earlier spring bloom (Adrian et al., 2006; Yang et al., 2016). According to

274    the ice record in Lake Erken (See Fig. S1, SI), in 2020, the lake was covered by very thin ice for only 5 days,

275    which is the shortest duration since observations were first recorded in 1954. The spring bloom in 2020 did occur

276    earlier than other years (See Fig. S3, SI), and both ML models which considered the timing of lake ice show fairly

277    good performance in predicting the timing and magnitude of this abnormally early spring bloom (Fig. 2, 5)

278    4.1.1 Performance of Hybrid PB ML models

279    One dimensional PB hydrodynamic models can accurately simulate both water temperature profiles, and other

280    hydrodynamic features in Lake Erken using the same forcing data that are commonly input to ML models. The

281  hybrid model structure tested here provides a richer set of input data leading to more accurate ML predictions of

282  algal *Chl* at little additional computational cost or data requirements. Using data from the hydrothermal PB model

283  allowed the seasonal deepening of the thermocline, variations in the surface mixing layer depth, and upwelling

284  events, represented by $W_n$, to be encoded into the ML algorithms. These factors can affect the underwater light

285  climate, the internal loading of phosphorus and the transport of resting cyanobacteria colonies from the

286  hypolimnion into the epilimnion favouring summer blooms of cyanobacteria (Pierson et al., 1992; Pettersson,

287  1998). The inclusion of these factors did increase the accuracy of the ML models, especially in the case of unusual

288  environmental conditions (e.g. spring of 2020, Fig. 2, 5) that did not frequently occur in the remaining

289  meteorological, hydrological and biogeochemical training data.

290  4.1.2 Prediction of bloom timing

291  For the purposes of water management, it may be most important to first predict the potential occurrence of a

292  bloom, and then once underway improve predictions of its magnitude. The best model performance in predicting

293  the timing of algal blooms, was obtained after adding hydrodynamic features derived from a PB model in

294  workflow 3, with TPR above 45% in detecting the onset of algal bloom during 2017-2020 and a modified accuracy

295  (Kappa) around 80 % indicated a moderate – strong level of prediction.

296  Based on our shuffling year tests of bloom timing, the GBR model showed relatively higher median TPRs than

297  LSTM model for sample intervals less than one month. However, in some training and testing year combinations,

298  TPRs are close to 0 % (Fig. 5), and CVs of the TPRs are highly variable, even at the original sample interval,

299  being over 30% for GBR and over 60% for LSTM, indicating that the correct detection of algal blooms in both

300  models are highly dependent on the years used to train the models. Thus, while the ML models can be better than

301  the PB models at predicting the onset of algal blooms, they still may not be good enough for operational

302  forecasting. The resulting variability provided a more accurate estimate of the model performance at each down-

303  sampled data interval and showed that increasing sample interval led to reduced performance for both ML models,

304  in terms of *MAE*, *RMSE*, and the CV of TPR. These tests also highlighted that the performance of both ML models,

305  especially LSTM, varied with the sampled history of events in the training period for evaluating a specific pattern

306  of change in the testing period. We suggest that testing strategies similar to the shuffle methods used in this study

307  are needed to accurately evaluate the expected accuracy of ML models when applied to any given site. The

308  estimated uncertainty in shuffling training year tests (Fig. 4) and shuffling training/testing year tests (Fig. 5) can

309  be used to better represent the uncertainty of ML derived forecasts.

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

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

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

### 5 Code availability

Model version 1.0 has been archived in Zenodo under DOI: 10.5281/zenodo.6534790, and is available at https://github.com/Shuqi-Lin/Erken_Algal_Bloom_Machine_Learning_Model.git.

### 6 Data availability

Lake Erken data are provided by Erken Laboratory, Uppsala University (https://www.ieg.uu.se/erken-laboratory/lake-monitoring-programme/, last access: May 2022). The dataset has been made possibly by the Swedish Infrastructure for Ecosystem Science (SITES), in this cast at Lake Erken.

337 **7 Supplement**

338 **8 Author contribution**

339 The concept of ML model workflow was designed by SL and DP. SL developed the ML model code and

340 performed the simulations. JM conducted the PB model simulations. SL wrote the manuscript with contributions

341 from DP and JM.

342 **9 Competing interests**

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

344 **10 Acknowledgement**

352 **References**

353 Adrian, R., Wilhelm, S., and Gerten, D.: Life-history traits of lake plankton species may govern their phenological response
354 to climate warming, Global Change Biology, 12, 652-661, 10.1111/j.1365-2486.2006.01125.x, 2006.
355 Baracchini, T., Wüest, A., and Bouffard, D.: Meteolakes: An operational online three-dimensional forecasting platform for
356 lake hydrodynamics, Water Research, 172, 115529, 10.1016/j.watres.2020.115529, 2020.
357 Brookes, J. D. and Carey, C. C.: Resilience to Blooms, Science, 334, 46-47, doi:10.1126/science.1207349, 2011.
358 Bruggeman, J. and Bolding, K.: A general framework for aquatic biogeochemical models, Environmental Modelling &
359 Software, 61, 249-265, https://doi.org/10.1016/j.envsoft.2014.04.002, 2014.
360 Burchard, H., Bolding, K., and Villarreal, M. R.: GOTM, a General Ocean Turbulence Model: Theory, Implementation and
361 Test Cases, European Commission. Joint Research Centre, Space Applications Institute, 103,
362 https://books.google.be/books/about/GOTM_a_General_Ocean_Turbulence_Model.html?id=zsJUHAAACAAJ&redir_esc=
363 y, 1999.
364 Burford, M. A., Carey, C. C., Hamilton, D. P., Huisman, J., Paerl, H. W., Wood, S. A., and Wulff, A.: Perspective:
365 Advancing the research agenda for improving understanding of cyanobacteria in a future of global change, Harmful Algae,
366 91, 101601, https://doi.org/10.1016/j.hal.2019.04.004, 2020.
367 Carey, C. C., Ibelings, B. W., Hoffmann, E. P., Hamilton, D. P., and Brookes, J. D.: Eco-physiological adaptations that
368 favour freshwater cyanobacteria in a changing climate, Water Research, 46, 1394-1407, 10.1016/j.watres.2011.12.016, 2012.
369 Elliott, J. A.: Is the future blue-green? A review of the current model predictions of how climate change could affect pelagic
370 freshwater cyanobacteria, Water Research, 46, 1364-1371, 10.1016/j.watres.2011.12.018, 2012.
371 Friedman, J. H.: Greedy Function Approximation: A Gradient Boosting Machine, The Annals of Statistics, 29, 1189-1232,
372 2001.
373 Hense, I. and Beckmann, A.: Towards a model of cyanobacteria life cycle—effects of growing and resting stages on bloom
374 formation of N2-fixing species, Ecological Modelling, 195, 205-218, https://doi.org/10.1016/j.ecolmodel.2005.11.018, 2006.
375 Hochreiter, S. and Schmidhuber, J.: Long Short-Term Memory, Neural Computation, 9, 1735-1780,
376 10.1162/neco.1997.9.8.1735, 1997.

377    Huisman, J., Codd, G. A., Paerl, H. W., Ibelings, B. W., Verspagen, J. M. H., and Visser, P. M.: Cyanobacterial blooms,
378    Nature Reviews Microbiology, 16, 471-483, 10.1038/s41579-018-0040-1, 2018.
379    Jimeno-Sáez, P., Senent-Aparicio, J., Cecilia, J. M., and Pérez-Sánchez, J.: Using Machine-Learning Algorithms for
380    Eutrophication Modeling: Case Study of Mar Menor Lagoon (Spain), International Journal of Environmental Research and
381    Public Health, 17, 1189, 2020.
382    Jöhnk, K. D., Brüggemann, R., Rücker, J., Luther, B., Simon, U., Nixdorf, B., and Wiedner, C.: Modelling life cycle and
383    population dynamics of Nostocales (cyanobacteria), Environmental Modelling & Software, 26, 669-677,
384    https://doi.org/10.1016/j.envsoft.2010.11.001, 2011.
385    Karlsson-Elfgren, I., Hyenstrand, P., and Riydin, E.: Pelagic growth and colony division of Gloeotrichia echinulata in Lake
386    Erken, Journal of Plankton Research, 27, 145-151, DOI 10.1093/plankt/fbh165, 2005.
387    Karlsson-Elfgren, I., Rengefors, K., and Gustafsson, S.: Factors regulating recruitment from the sediment to the water
388    column in the bloom-forming cyanobacterium Gloeotrichia echinulata, Freshwater Biology, 49, 265-273, DOI
389    10.1111/j.1365-2427.2004.01182.x, 2004.
390    Marcé, R., George, G., Buscarinu, P., Deidda, M., Dunalska, J., de Eyto, E., Flaim, G., Grossart, H.-P., Istvanovics, V.,
391    Lenhardt, M., Moreno-Ostos, E., Obrador, B., Ostrovsky, I., Pierson, D. C., Potužák, J., Poikane, S., Rinke, K., Rodríguez-
392    Mozaz, S., Staehr, P. A., Šumberová, K., Waajen, G., Weyhenmeyer, G. A., Weathers, K. C., Zion, M., Ibelings, B. W., and
393    Jennings, E.: Automatic High Frequency Monitoring for Improved Lake and Reservoir Management, Environmental Science
394    & Technology, 50, 10780-10794, 10.1021/acs.est.6b01604, 2016.
395    McHugh, M. L.: Interrater reliability: the kappa statistic, Biochemia medica, 22, 276-282, 2012.
396    Mesman, J. P., Ayala, A. I., Goyette, S., Kasparian, J., Marcé, R., Markensten, H., Stelzer, J. A. A., Thayne, M. W., Thomas,
397    M. K., Pierson, D. C., and Ibelings, B. W.: Drivers of phytoplankton responses to summer wind events in a stratified lake: A
398    modeling study, Limnology and Oceanography, 67, 856-873, https://doi.org/10.1002/lno.12040, 2022.
399    Moras, S., Ayala, A. I., and Pierson, D. C.: Historical modelling of changes in Lake Erken thermal conditions, Hydrology
400    and Earth System Sciences, 23, 5001-5016, 2019.
401    Nelson, N. G., Muñoz-Carpena, R., Phlips, E. J., Kaplan, D., Sucsy, P., and Hendrickson, J.: Revealing Biotic and Abiotic
402    Controls of Harmful Algal Blooms in a Shallow Subtropical Lake through Statistical Machine Learning, Environmental
403    Science & Technology, 52, 3527-3535, 10.1021/acs.est.7b05884, 2018.
404    Paerl, H. W.: Nuisance phytoplankton blooms in coastal, estuarine, and inland waters1, Limnology and Oceanography, 33,
405    823-843, 10.4319/lo.1988.33.4part2.0823, 1988.
406    Paerl, H. W. and Huisman, J.: Blooms Like It Hot, Science, 320, 57-58, doi:10.1126/science.1155398, 2008.
407    Persson, I. and Jones, I. D.: The effect of water colour on lake hydrodynamics: a modelling study, Freshwater Biology, 53,
408    2345-2355, https://doi.org/10.1111/j.1365-2427.2008.02049.x, 2008.
409    Pettersson, K.: The Availability of Phosphorus and the Species Composition of the Spring Phytoplankton in Lake Erken,
410    Internationale Revue der gesamten Hydrobiologie und Hydrographie, 70, 527-546, 10.1002/iroh.19850700407, 1985.
411    Pettersson, K.: Mechanisms for internal loading of phosphorus in lakes, Hydrobiologia, 373, 21-25,
412    10.1023/A:1017011420035, 1998.
413    Pettersson, K., Grust, K., Weyhenmeyer, G., and Blenckner, T.: Seasonality of chlorophyll and nutrients in Lake Erken –
414    effects of weather conditions, Hydrobiologia, 506, 75-81, 10.1023/B:HYDR.0000008582.61851.76, 2003.
415    Pierson, D. C., Pettersson, K., and Istvanovics, V.: Temporal changes in biomass specific photosynthesis during the summer:
416    regulation by environmental factors and the importance of phytoplankton succession, Hydrobiologia, 243, 119-135,
417    10.1007/BF00007027, 1992.
418    Read, J. S., Hamilton, D. P., Jones, I. D., Muraoka, K., Winslow, L. A., Kroiss, R., Wu, C. H., and Gaiser, E.: Derivation of
419    lake mixing and stratification indices from high-resolution lake buoy data, Environmental Modelling & Software, 26, 1325-
420    1336, 10.1016/j.envsoft.2011.05.006, 2011.
421    Recknagel, F., Fukushima, T., Hanazato, T., Takamura, N., and Wilson, H.: Modelling and prediction of phyto- and
422    zooplankton dynamics in Lake Kasumigaura by artificial neural networks, Lakes & Reservoirs: Science, Policy and
423    Management for Sustainable Use, 3, 123-133, 10.1111/j.1440-1770.1998.tb00039.x, 1998.
424    Reichwaldt, E. S. and Ghadouani, A.: Effects of rainfall patterns on toxic cyanobacterial blooms in a changing climate:
425    Between simplistic scenarios and complex dynamics, Water Research, 46, 1372-1393, 10.1016/j.watres.2011.11.052, 2012.
426    Richardson, J., Miller, C., Maberly, S. C., Taylor, P., Globevnik, L., Hunter, P., Jeppesen, E., Mischke, U., Moe, S. J.,
427    Pasztaleniec, A., Søndergaard, M., and Carvalho, L.: Effects of multiple stressors on cyanobacteria abundance vary with lake
428    type, Global Change Biology, 24, 5044-5055, 10.1111/gcb.14396, 2018.
429    Rousso, B. Z., Bertone, E., Stewart, R., and Hamilton, D. P.: A systematic literature review of forecasting and predictive
430    models for cyanobacteria blooms in freshwater lakes, Water Research, 182, 115959, 10.1016/j.watres.2020.115959, 2020.
431    Stanley, F. K. T., Irvine, J. L., Jacques, W. R., Salgia, S. R., Innes, D. G., Winquist, B. D., Torr, D., Brenner, D. R., and
432    Goodarzi, A. A.: Radon exposure is rising steadily within the modern North American residential environment, and is
433    increasingly uniform across seasons, Scientific Reports, 9, 18472, 10.1038/s41598-019-54891-8, 2019.
434    Watson, S. B., Miller, C., Arhonditsis, G., Boyer, G. L., Carmichael, W., Charlton, M. N., Confesor, R., Depew, D. C.,
435    Höök, T. O., Ludsin, S. A., Matisoff, G., McElmurry, S. P., Murray, M. W., Peter Richards, R., Rao, Y. R., Steffen, M. M.,
436    and Wilhelm, S. W.: The re-eutrophication of Lake Erie: Harmful algal blooms and hypoxia, Harmful Algae, 56, 44-66,
437    https://doi.org/10.1016/j.hal.2016.04.010, 2016.
438    Wei, B., Sugiura, N., and Maekawa, T.: Use of artificial neural network in the prediction of algal blooms, Water Research,
439    35, 2022-2028, 10.1016/S0043-1354(00)00464-4, 2001.
440    Xiao, X., He, J., Huang, H., Miller, T. R., Christakos, G., Reichwaldt, E. S., Ghadouani, A., Lin, S., Xu, X., and Shi, J.: A
441    novel single-parameter approach for forecasting algal blooms, Water Research, 108, 222-231, 10.1016/j.watres.2016.10.076,
442    2017.

443     Yang, Y., Stenger-Kovács, C., Padisák, J., and Pettersson, K.: Effects of winter severity on spring phytoplankton
444     development in a temperate lake (Lake Erken, Sweden), Hydrobiologia, 780, 47-57, 10.1007/s10750-016-2777-8, 2016.
445
446