# The predictability of a lake phytoplankton community, from hours to years

Mridul K. Thomas[1,2*], Simone Fontana[1, 3], Marta Reyes[1], Michael Kehoe[4], and Francesco Pomati[1]

*Corresponding author. Address: Centre for Ocean Life, DTU Aqua, Technical University of Denmark, Kemitorvet, 2800 Kongens Lyngby, Denmark. Phone: +45 3588 3412

[1] Dept. of Aquatic Ecology, Eawag: Swiss Federal Institute of Aquatic Science and Technology, Überlandstrasse 133, 8600 Dübendorf, Switzerland
[2] Centre for Ocean Life, DTU Aqua, Technical University of Denmark, Kemitorvet, 2800 Kongens Lyngby, Denmark.
[3] Biodiversity and Conservation Biology Research Unit, Conservation Biology Group, Swiss Federal Institute of Forest, Snow and Landscape Research, Zürcherstrasse 111, 8903 Birmensdorf, Switzerland
[4] Global Institute for Water Security and School of Environment and Sustainability, University of Saskatchewan

**Author email IDs:**
Mridul K. Thomas: mrit@dtu.dk
Simone Fontana: simone.fontana@wsl.ch
Marta Reyes: marta.reyes@eawag.ch
Michael Kehoe: kehoe.michael@eawag.ch
Francesco Pomati: francesco.pomati@eawag.ch

**Key words**: prediction, forecast, machine learning, phytoplankton, time series, environmental monitoring, cyanobacteria

**Statement of authorship**: MKT & FP conceived the study. MKT, SF, MR & FP collected the data. MKT analysed the data. MKT wrote the manuscript with substantial input from FP, SF, MK & MR.

**Abstract**

Forecasting anthropogenic changes to ecological communities is one of the central challenges in ecology. However, nonlinear dependencies, biotic interactions and data limitations have limited our ability to assess how predictable communities are. Here we used a machine learning approach and environmental monitoring data (biological, physical and chemical) to assess the predictability of phytoplankton cell density in one lake across an unprecedented range of time scales. Communities were highly predictable over hours to months: model $R^2$ decreased from 0.89 at 4 hours to 0.75 at 1 month, and in a long-term dataset lacking fine spatial resolution, from 0.46 at 1 month to 0.32 at 10 years. When cyanobacterial and eukaryotic algal cell density were examined separately, model-inferred environmental growth dependencies matched laboratory studies, and suggested novel trade-offs governing their competition. High-frequency monitoring and machine learning can help elucidate the mechanisms underlying ecological dynamics and set prediction targets for process-based models.

## Introduction

Forecasting how environmental change will alter communities and ecosystems is perhaps the most important task facing ecologists today, and a tremendous challenge to our ecological understanding (Mouquet et al. 2015, Petchey et al. 2015, Houlahan et al. 2016). Nonlinear relationships (such as between temperature and most biological processes), stochasticity and sensitive dependence on initial conditions are sources of uncertainty that community ecology shares with other predictive disciplines, such as climate science. However, ecology additionally has to grapple with biotic interactions in complex food webs, evolutionary change, and a paucity of data with which to assess predictive power and refine models (Magurran et al. 2010). Therefore, with few exceptions (notably in disease ecology, e.g. Axelsen et al. 2015), we do not know how predictable ecological communities are (i.e. how strong the association between present and future system states is, a proxy for forecast ability). Quantifying this would allow us to understand the time scale over which we can provide actionable input for management and legislative decision-making, recently termed the *ecological forecast horizon* (Petchey et al. 2015).

To make accurate long-term forecasts, ecology needs to develop process-based forecasts akin to those prevalent in climate science. Correlational approaches based on present conditions and abundances are likely to make inaccurate forecasts over decadal time-scales because patterns of environmental covariation will change in the future (Williams et al. 2007). Process-based models avoid this problem but are a challenge to design because of their complexity. This arises from of a lack of knowledge of the functional forms (or shape) relating population/community change to environmental factors, and a lack of data with which to parameterise them (Kremer et al. 2016). The scale of this challenge is highlighted by recent work showing that complex, highly nonlinear interactions between abiotic factors are a regular feature of physiological and ecological processes (Zhu et al. 2016, Zimmer et al. 2016, Edwards et al. 2016, Thomas et al.

3

85    2017). Designing process-based models using a traditional approach may require extensive

86    experimental work examining high-dimensional interactions. High-throughput screening

87    technologies are helping to address this problem. But in many cases, a traditional experimental

88    approach to understanding interactions (i.e. through multidimensional factorial experiments) may

89    not be realistic given present funding and experimental constraints.

90    Machine learning (ML) algorithms offer us an alternative path towards the creation of these

91    process-based models. When presented with complex environmental datasets, ML allows us to

92    avoid the most important constraints inherent in traditional statistical approaches (*a priori*

93    specification of functional forms, interactions and error distributions). Despite relying on

94    underlying correlations, ML algorithms can improve substantially on traditional correlative

95    analyses (Rivero-Calle et al. 2015, Kehoe et al. 2012, 2016). They can be used on complex

96    datasets to assess associations (Rivero-Calle et al. 2015) and quantify predictability in the

97    absence of the knowledge needed for a process-based model (Ewers et al. 2017). Even more

98    importantly, they can be used to *infer* the functional forms and interactions needed to develop

99    process-based models. This approach will require large datasets, but as ecology enters the 'big

100   data' era, acquiring this is becoming feasible for many systems. As the cost of data acquisition

101   continues to decrease, ML may prove a more efficient approach (relative to high-dimensional

102   factorial experiments) to assessing community predictability and understanding the drivers of

103   complex ecological dynamics.

104   Natural communities of microbes such as phytoplankton are vital parts of most biogeochemical

105   cycles and food webs (Field et al. 1998, Falkowski et al. 1998), and so assessing their

106   predictability is especially important. Phytoplankton have generation times on the order of a day,

107   and respond extremely rapidly to environmental change: shifts on time-scales of minutes to

108   hours are sufficient to elicit physiological and ecological changes (Goldman & Glibert 1982,

4

109  Demers et al. 1991, Hemme et al. 2014). Despite this sensitivity to environmental conditions, we

110  do not know the time-scales over which phytoplankton community dynamics may be predicted.

111  Historically, most plankton monitoring campaigns have measured the community at coarse time-

112  scales of once to twice a month (Jochimsen et al. 2012), amounting to tens of generations.

113  These efforts have helped us understand broad changes driven by eutrophication and

114  environmental warming (Pomati et al. 2012, Jochimsen et al. 2012), helping to make the case for

115  policies limiting further changes. However, with rare exceptions (notably Hunter-Cevera et al.

116  2014, 2016), plankton monitoring efforts have not captured data needed to accurately assess

117  community predictability across time scales. High-frequency monitoring campaigns that sample

118  communities and environmental drivers on sub-daily time-scales can partly address this (Pomati

119  et al. 2011, Merel 2013, Pomati et al. 2013, Hunter-Cevera et al. 2014, 2016), filling in pieces of

120  the picture that coarser long-term datasets have hinted at. They also provide us with the quantity

121  of data needed to profitably employ ML tools.

122  We quantified the predictability of phytoplankton cell density over time scales ranging from 4

123  hours to 10 years, or approximately $10^{-1}$ to $10^3$ generations. Cell density, or the abundance of

124  phytoplankton cells per unit volume, is the most important parameter characterising the

125  phytoplankton community. It is a strong proxy for phytoplankton biomass (including in our

126  system, Fig. S1) and for primary productivity, an important ecosystem property. We

127  characterised predictability of cell density in Greifensee, a meso-eutrophic peri-alpine lake in

128  Switzerland. The Greifensee plankton community, chemistry and physics have been monitored

129  for >30 years, during which it has seen dramatic changes in biology as a result of eutrophication

130  and re-oligotrophication (Bürgi et al. 2003). We make use of two complementary datasets

131  examining the Greifensee phytoplankton community: (1) High-frequency data from monitoring

132  campaigns carried out in the summer and autumn of 2014 and 2015. Cell density was measured

133  at 6 different depths (Fig. 1) using scanning flow cytometry (SFCM), and environmental data

134    were also collected (Table S1). (2) A long-term time series created from monthly measurements

135    of depth-integrated phytoplankton measurements (Fig. 1), as well as associated environmental

136    factors (Table S2). Although the datasets differ in methodology and size, they measure

137    substantially similar biological, physical and chemical factors (Tables S1, S2). Given the length

138    of the two datasets and their sampling frequency and location, we are able to directly compare

139    predictability in the two datasets at a time lag of 1 month.

140    In addition to community density, ecology aims to predict the dynamics of functional groups and

141    taxa. Especially because toxic cyanobacterial blooms are a major health concern (Chorus &

142    Bartram 1999, Paerl & Huisman 2009, Paerl et al. 2011, Merel et al. 2013), we also assessed

143    the predictability of cyanobacterial cell density over these time scales, and the drivers of

144    competitive dynamics between cyanobacteria and eukaryotic phytoplankton (Fig. 2). Eukaryotic

145    densities have remained relatively stable in Greifensee since the 1980s, while average

146    cyanobacterial densities first increased 100-fold during eutrophication and then decreased by a

147    similar amount over this time period as a result of re-oligotrophication (Fig. 2). Understanding the

148    drivers of growth and competition between these two broad phytoplankton groups can help us

149    refine process-based models of water quality, with important implications for the management of

150    aquatic ecosystem services.

151    **Methods**

152    **I. Overview**

153    Greifensee is a peri-alpine lake in Switzerland (47.35°N, 8.68°E) with a documented history of

154    eutrophication and re-oligotrophication (Bürgi et al. 2003). The lake is currently meso-eutrophic,

155    32 m deep at its deepest point and just over 20 m deep at the sampling locations used for both

156    datasets.

157     We use two datasets in this study: a high-frequency dataset consisting of measurements every 4

158     hours during the summer and autumn of 2014 and 2015 using the automated monitoring station

159     Aquaprobe (Pomati et al. 2011), and a long-term dataset consisting of monthly measurements

160     from March 1984 to June 2016. In both cases, important environmental data (both abiotic and

161     biotic) was collected simultaneously near the middle of the lake, allowing similar analyses to be

162     conducted and thereby enabling comparisons. However, the datasets differ in important ways:

163     i) The high-frequency dataset involved measurements by SFCM, while the long-term dataset

164     involved microscopy measurements. Therefore, sampling effort and density assessment

165     methods differ.

166     ii) The high-frequency dataset consists of measurements at six specific depths (1.0, 2.5, 4.0, 5.5,

167     7.0 and 8.5m). Abiotic environmental data were also estimated at the same depth as the

168     collected sample. In contrast, the long-term dataset consists of integrated phytoplankton

169     measurements across the top 20m of the lake. Abiotic measurements were not integrated, but

170     collected at specific depths (except for light, which is a surface estimate), and so we calculated

171     the maximum and minimum value of each abiotic factor in the top 20m for use as our predictors.

172     iii) The high-frequency dataset consists of 7161 measurements, while the long-term dataset

173     contains 383 measurements.

174     **II. High-frequency dataset generation**

175     *1. Scanning flow cytometry description*: We used a scanning flow cytometer, the CytoSense

176     (http://www.cytobuoy.com), to quantify the density of the total phytoplankton community as well

177     as its cyanobacterial and eukaryotic algal fractions (estimated densities are strongly correlated

178     with estimates from microscopy, Fig. S2). The CytoSense characterizes the scattering and

179     pigment fluorescence of individual phytoplankton cells. It measures cells and colonies across a

180     large proportion of the phytoplankton length range, between approximately 2 µm and 1 mm in

181    length. Particles that enter the system cross two coherent 15mW solid-state lasers. The

182    instrument's laser and sensor wavelengths are designed to target the fluorescence signals

183    primarily from chlorophyll-a and phycocyanin, but also capture signals from phycoerythrin and

184    carotenoids. We used two different instruments 2014 and 2015, with small differences in

185    configuration. Instrument settings and data processing steps may be found in the supplementary

186    information.

187    *2. SFCM field sampling procedure*: All samples were collected from a floating platform

188    Aquaprobe (Pomati et al. 2011) near the middle of the lake (47.3663°N, 8.665°E). Every four

189    hours, water was sampled automatically at each of the 6 depths (as described in Pomati et al.

190    2011). Water samples were pumped into a 150 mL sampling chamber at the surface through a

191    tube with a 0.6-cm diameter opening. The sampling chamber was flushed with water from the

192    sampling depth three to five times over 2 minutes before the CytoSense collected a subsample

193    of up to 500 µL for measurement.

194    *3. Environmental factors:* The full list of environmental parameters is found in Table S1. We

195    measured temperature, conductivity and irradiance at all depths. We also collected weekly

196    depth-specific samples for dissolved nutrient (nitrate, phosphate, ammonium) concentration

197    estimation, and integrated measurements of size-fractionated zooplankton. We also monitored

198    meteorological factors including wind speed and rainfall, and made use of additional data

199    provided by the Office of Waste, Water, Energy and Air (AWEL) of Canton Zürich on water inflow

200    (including flow rate, temperature and nutrient concentrations) into the lake. Detail of sampling

201    procedures, instruments used and measurement methodology may be found in the

202    supplementary information.

203

204

**III. Long-term dataset**

205

206    *1. Field sampling procedure*: Approximately every month, water samples were collected for

207    physical, chemical and biological measurements near the centre of the lake (47.3525°N,

208    8.6748°E), approximately 1.5 km from the floating platform used for the high-frequency

209    measurements. For microscopic counts of the phytoplankton community, an integrated water

210    sample was collected over the upper 20 m of the water column with a Schröder sampler (Bürgi

211    et al. 2003).

212    *2. Environmental factors*: The full list of parameters is found in Table S2. To measure chemical

213    and physical parameters, water samples were collected every 2.5 m over the whole water

214    column, at the same location and on the same dates, and were analysed using standard

215    limnological methods (Rice et al. 2012). Integrated zooplankton samples were collected over the

216    upper 20 m of the water column. We also made use of a publicly available surface irradiance

217    dataset (Schulz et al. 2008, Müller et al. 2015) to estimate the monthly-averaged irradiance at

218    the water surface based on interpolated estimates from a location approximately 2 km from the

219    sampling location (47.35°N 8.65°E). More details about sampling procedures, instruments used

220    and measurement methodology may be found in Bürgi et al. (2003).

221    *3. Data processing*: For every time point, we calculated the maximum and minimum value of

222    every depth-specific parameter (such as phosphate concentration) across the entire water

223    column, and used these for subsequent analyses. Additionally, samples were not collected on

224    the same day every month and, in rare cases, more than one sample was collected in a month.

225    We therefore aggregated measurements by rounding to the nearest month and then averaged

226    duplicate values.

227    **IV. Machine learning analyses**

9

228    *1. Random forests overview*: Random forests (RFs) are a robust ML tool comprising ensembles

229    of regression trees (or classification trees) (Breiman 1999). In each regression 'tree' within the

230    random 'forest', a randomly selected subset of the data is recursively partitioned based on the

231    most strongly associated predictor. At each node, a random subset of the total number of

232    predictors is considered for partitioning. The final tree 'prediction' for new data is given by the

233    average value of the data within each branch of the tree. By aggregating predications across

234    trees, RFs are able to reproduce arbitrarily complex shapes patterns without *a priori* functional

235    form specification.

236    We took advantage of three features that make RFs a flexible and useful tool for examining

237    ecological systems: *permutation importance*, easy quantification of *partial effects* of individual

238    predictors, and *out-of-bag prediction*.

239     (i) The importance of each predictor in a RF is assessed by permuting the predictor across all

240    trees in the forest and quantifying the resulting change in the forest's error rate. More important

241    predictors lead to a greater increase in error when permuted.

242    (ii) The partial effect of any single predictor on the dependent variable can also be quantified,

243    allowing us to examine the functional form of the relationship (which may be arbitrarily nonlinear,

244    though non-bifurcating).

245    (iii) Out-of-bag (OOB) prediction allows us to make accurate estimates of error rate and

246    goodness of fit (model $R^2$) via a process akin to cross validation (Breiman 1999). Each data

247    point is present in the training data of only a subset of all 'trees' that comprise the 'forest'.

248    Therefore, the value of every point may be predicted using the trees that have not been trained

249    with it. The 'OOB prediction error', or mean difference between the OOB predictions and the true

250    value of all points in the dataset (see Fig. S3, S4 for examples using our data) can be used to

251    quantify the RF's predictive ability through a pseudo-$R^2$:

10

252
$$R^2 = 1 - \frac{MSE}{var(y)}$$

253     where *MSE* is the mean squared error of the OOB predictions when compared to the true

254     values, and *var(y)* is the variance in the dependent variable. As in the case of a standard $R^2$, a

255     pseudo-$R^2$ has an upper bound of 1, indicating perfect model performance. However, note that

256     unlike a standard $R^2$, there is no lower bound. It is possible for the pseudo-$R^2$ to be negative, if

257     *MSE > var(y)*. This may be interpreted as saying that the model prediction is worse than the

258     mean value of the dependent variable in the entire dataset. In our analyses, we saw low,

259     negative values of pseudo-$R^2$ in a small number of cases; we rounded these values to zero to

260     avoid confusion, while noting this in the figure captions.

261     2. *Data pre-processing*: In our analyses, we omitted: i) all entries cases where the dependent

262     variable was missing, and ii) the predictors *sampling depth* and *sampling time*. We omitted the

263     latter in order to accurately estimate the effects of predictors that covary with depth and time on

264     cell density. In other words, we believe that gradients in light, temperature and nutrients should

265     characterise most of the relevant information contained within depth and time.

266     3. *Quantifying predictability:* We quantified the *predictability* of log cell density of the total

267     phytoplankton community, and the cyanobacterial and eukaryotic fractions, using the pseudo-$R^2$

268     calculated based on the OOB predictions of the fitted model (see details above). We estimated

269     predictability at time lags ranging from 4 hours to 1 month in the high-frequency dataset, and

270     from 1 month to 10 years in the long-term dataset. For every time lag, we fit two models,

271     predicting log cell density using: 1) only log cell density at the specified time lag, and 2) both log

272     cell density and environmental parameters at the specified time lag. E.g. our simplest model

273     considering a time lag of four hours predicted log cell density at all time points using only log cell

274     density from the measurement four hours previously.

11

275  *4. Predictor importance*: We assessed the importance of predictors at all time lags using the

276  change in model error rate when the predictor values were permuted.

277  *5. Partial effects of environment on growth*: We quantified the model-inferred effects of

278  environmental factors on cyanobacteria and eukaryotes. Instead of examining the effects of

279  these predictors on log density, we instead examined how they influence the population growth

280  rate (i.e. specific growth rate, day$^{-1}$, the rate of change in density between successive time

281  points). We did this to facilitate comparison between the partial effects in our field dataset and

282  extensive prior laboratory findings for the same predictors. However, the functional forms

283  remained highly similar to the model explaining log density.

284  We focussed on two factors that are known to influence phytoplankton growth (Litchman &

285  Klausmeier 2008) and were identified as important in our analyses: light and temperature.

286  Because laboratory studies typically measure the effects of environmental factors on population

287  growth rate *per day*, we multiplied the estimates of growth rate over four hours by 6 to express

288  them in the same units. We then fitted RFs to these population growth rates using environmental

289  parameters at a 4-hour lag and estimated their partial effects. Note that we omitted log density

290  as a predictor, but model structure was otherwise identical to those previously described.

291  Though we were also interested in the effects of dissolved nitrate, phosphate and N:P ratio, we

292  had less well-resolved data for these predictors that limited the power of analyses relating to

293  these factors.

294  *5. Model fitting and settings*: All analyses were done in the *R* statistical environment v3.3.3 (R

295  Core Team 2017) using the package *randomforestSRC* (Ishwaran & Kogalur 2007, Ishwaran &

296  Kogalur 2017). We used 2000 trees for every forest, and set the number of predictors to be

297  considered at each node to be one-third of the total number of predictors. Missing data among

298  the predictors were imputed for the purpose of model fitting, but imputed values were not used

299  for predictor importance assessment (Ishwaran & Kogalur 2007, Ishwaran & Kogalur 2017).

## Results

301  Phytoplankton cell density was highly predictable on time scales of hours to months. In our high-

302  frequency dataset, pseudo-$R^2$ of the RF models trained with cell density and environmental data

303  decreased from 0.89 at a 4 hour lag to 0.74 at a lag of 1 month (Fig. 3). The model using only

304  cell density as a predictor had a lower $R^2$ at all time lags. As time lag increased, including

305  environmental data led to larger improvements in predictability: the difference in $R^2$ between the

306  two models was 0.03 at a 4 hour lag, and ten times higher (0.30) at a lag of 1 month (Fig. 3). In

307  the long-term dataset, $R^2$ of the model trained with cell density and environmental data

308  decreased from 0.46 at a time lag of 1 month to 0.35 at 6 months, after which it remained

309  relatively stable (Fig. 3). $R^2$ in the density-only model was lower at all time lags.

310  Aside from cell density, which was the strongest predictor at all time lags in our high-frequency

311  dataset, the most important predictors were light, temperature and thermocline depth, itself an

312  indirect effect of temperature (Fig. 4). Light and temperature were also most important in our

313  long-term dataset on time scales of months (Fig. 4). At time-scales of years, dissolved

314  phosphorus and zooplankton density become more important.

315  Cyanobacteria were more predictable than eukaryotes at all time scales, in both high-frequency

316  and long-term datasets (Fig. 3). Model $R^2$ for cyanobacteria was consistently higher that than for

317  eukaryotes by approximately 5-20 percentage points (Fig. 3), in both types of models (cell

318  density only and cell density with environmental data).

319  To motivate the development of process-based models of phytoplankton competition, we also

320  examined the partial effects of environmental factors on the growth rate of cyanobacteria and

321  eukaryotic algae (Fig. 5). Temperature and light, the strongest predictors in our dataset, showed

13

322    biologically realistic nonlinear patterns. Additionally, in both cases, each group dominated a

323    region of parameter space, suggesting the presence of trade-offs in performance.

### Discussion

325    Assessing the predictability of natural communities is crucial if we are to develop forecasts of

326    how ecosystems will be altered by anthropogenic environmental change (Petchey et al. 2015,

327    Mouquet et al. 2015, Houlahan et al. 2016). However, our ability to predict community dynamics

328    has been limited by our understanding of environmental dependencies and biotic interactions

329    (McGill et al. 2006). Our results suggest that lake phytoplankton communities are highly

330    predictable over time scales of hours to months, approximately $10^{-1}$ to $10^2$ generations, and

331    possibly longer (Figs. 3, S5). Our approach quantifies the decline in predictability with increasing

332    time lag, identifies the predictors that contribute to predictive power, and points towards realistic

333    trade-offs and parameterisations through the examination of partial effects. Together, these can

334    inform the development of process-based models, set targets for their forecasts to achieve, and

335    identify a forecast horizon for adaptive management strategies. This is especially true in the

336    case of cyanobacteria, which are a threat to human health and aquatic ecosystem services

337    because of toxin production, and are believed to be hard to forecast (Chorus & Bartram 1999,

338    Paerl & Huisman 2009, Paerl et al. 2011, Merel et al. 2013). We find cyanobacterial densities to

339    be consistently more predictable than those of eukaryotes (Fig. 3).

340    As our understanding of ecological processes improves, the limits to predictability of ecological

341    systems will be determined by more fundamental constraints such as stochasticity, and sensitive

342    dependence on initial conditions. Despite these forces, we find strong, ecologically important

343    environmental forcing in a natural system across a range of time scales (Figs. 3-5).

344    Consequently, we believe that process-based models are very likely to provide us with useful

345    predictions over medium-to-long time scales. In other words, we believe that despite the

346    complexity of phytoplankton communities, the ecological forecast horizon (Petchey et al. 2015)

14

347    is sufficiently distant for ecologists to provide useful input into adaptive management strategies.

348    Note that we do not quantify a specific horizon here because this requires the specification of a

349    (arbitrary) forecast threshold; if desired, readers may choose these thresholds for themselves

350    and identify the resulting forecast horizon from Fig. 3.

351    Light and temperature were strongly predictive of phytoplankton dynamics across time scales

352    (Fig. 4), consistent with existing ecological understanding (Litchman & Klausmeier 2008). We

353    also found that zooplankton density and dissolved phosphorus concentrations become highly

354    predictive on time scales longer than a year, consistent with an ongoing, multidecadal decrease

355    in phosphorus and biomass in Greifensee (Buergi et al. 2003). This identification of variables

356    that are known to play a major role in phytoplankton ecology strengthens our confidence in the

357    relationships underlying our metric of predictability. However, we note that predictor importance

358    – while a useful tool – is sensitive to missing data patterns. Our estimates therefore understate

359    the importance of two major groups of predictors in our high frequency data: nutrients and

360    zooplankton density. Unlike most other predictors that were measured every four hours, these

361    were measured weekly in 2014 and twice a week in 2015 (Table S1). To partially correct for this

362    difference, we also assessed the relative importance of all predictors when these were

363    interpolated using generalised additive models (GAMs) (Fig. S6). Models with interpolated

364    nutrients and zooplankton predictors had marginally higher $R^2$ values and these predictors rose

365    considerably in importance, especially dissolved nitrogen. We believe that these results are

366    noteworthy, but choose not to focus on them here because we are unable to validate the

367    interpolated estimates.

368    Importantly, the predictive power of environmental factors in our models arises from nonlinear

369    dependencies that are consistent with causal relationships established through lab studies (Fig.

370    5; Litchman & Klausmeier 2008). Light, one of the most important predictors, has a partial effect

371    on growth that is a saturating function for cyanobacteria and a right-skewed unimodal function

372    for eukaryotes (Fig. 5); these are the only shapes consistent with laboratory measurements of

373    light-dependent growth (Eilers & Peeters 1998, Edwards et al. 2015). The partial effect of

374    temperature is an increasing function and possibly a left-skewed unimodal curve, consistent with

375    prior eco-physiological findings, including in phytoplankton (Kingsolver 2009, Thomas et al.

376    2012, Thomas et al. 2016). This concordance between controlled lab studies and ML-derived

377    field patterns increases our confidence in the suitability of this ML approach, and suggests that

378    the relationships we have uncovered are likely to be useful in guiding process-based model

379    creation. Furthermore, it suggests that ML approaches may be used to discover novel ecological

380    patterns. This is particularly important in the case of interactions between factors, which

381    presently require labour-intensive and expensive multifactorial experiments to understand.

382    The partial effects that we show here (Fig. 5) point towards trade-offs that could enable the co-

383    existence of cyanobacteria and eukaryotes. Cyanobacteria appear to benefit from high light

384    intensity and high temperature, while eukaryotes have a growth advantage in the converse

385    conditions. Therefore, temporal heterogeneity in one or both of these dimensions could allow for

386    the maintenance of both these groups (Chesson 2000). Cyanobacteria do possess higher

387    optimal temperatures for growth than eukaryotic phytoplankton at temperate latitudes (Thomas

388    et al. 2016), consistent with the temperature-dependence we see (Fig. 5). The apparent trade-off

389    between growth at high and low light intensities was not seen in a synthesis of lab-measured

390    light traits (Schwaderer et al. 2011), but at present, measurements are available only from a

391    small number of species and may be influenced by interactions with other factors. Laboratory

392    data on a broader range of species and under a greater range of conditions will be needed to

393    resolve this discrepancy. If true, the pattern we observe in the field also suggests an explanation

394    underlying the formation of surface scums by cyanobacteria through buoyancy regulation (Paerl

395    et al. 2011, Carey et al. 2012). Scum formation - important due to the negative impact on lake

396    ecosystem services - is consistent with a cyanobacterial benefit from higher irradiance. In

397    contrast, eukaryotes appear to have a lower optimal irradiance and might experience photo-

16

398  degradation from surface growth. These observations offer an example of the insights that may

399  be gained through a combination of high-frequency monitoring and machine learning.

400  Our models may understate the long-term predictability of the phytoplankton community. The

401  difference in predictability between high-frequency and long-term datasets at a time lag of 1

402  month suggests that if a similar methodology was followed in the long-term dataset, reasonably

403  high $R^2$ values may have been obtained over time scales of years, not just months. This

404  difference is driven by several factors: 1) our high-frequency dataset includes measurements of

405  both phytoplankton and environmental factors at specific depths, as opposed to integrated

406  values across the water column as in the long-term dataset, 2) the high-frequency dataset has

407  more than an order of magnitude more data points (7161 vs. 383) with which to train the

408  machine learning algorithm, and 3) the long-term dataset explores a far greater range of

409  parameter space in temperature, nutrient concentration, zooplankton density and unmeasured

410  factors. Of the three, we believe depth-specific sampling may be the major factor, as the

411  difference in model $R^2$ at a lag of one month is <15% in the case of the cell density-only models,

412  and 30% in the models with both cell density and environmental factors (Fig. 3). However, we

413  also note that in more complex systems where migration is a larger factor – such as coastal and

414  open-ocean communities – predictability may be lower unless physical circulation patterns are

415  highly predictable as well.

416  It is important to note that although pseudo-$R^2$ provides estimates of predictability that are robust

417  (Breiman 1999), we have not assessed a true *forecast*, in which error is allowed to compound

418  through time. This approach can in principle be used to make a forecast, but we chose not to do

419  so because of large changes in environmental conditions towards the end of the 2014 and 2015

420  monitoring seasons. Attempting to forecast would require us to predict in conditions well outside

421  those that the model was trained on, where it will inevitably perform poorly. Despite this

422  limitation, we believe that out-of-bag error is a useful proxy for forecast error: the realistic

17

423  environmental dependencies (Fig. 5) highlight that we are uncovering the mechanisms

424  underpinning ecological dynamics. In the future, a broader sampling of parameter space

425  (through a year-round monitoring campaign) should allow us to make and test true forecast skill.

426  We have shown that high-frequency environmental monitoring and machine learning

427  approaches can be usefully employed to uncover patterns in complex ecological communities, to

428  assess the predictability of these communities, and to uncover dependencies that can then be

429  incorporated into process-based models of communities and ecosystems. This can help us

430  address fundamental questions in ecology: What are the drivers of ecological processes and

431  how does this change through time? How large of an effect does environmental and

432  demographic stochasticity have on communities? What are the dominant trade-offs that maintain

433  diversity in natural systems and how do they operate in dynamic environments? But perhaps

434  more importantly, it can allow us to improve our forecasts of ecological systems, fulfilling a

435  fundamental obligation that ecology owes to society.

## Acknowledgements

445

446

18

447 **References**

448 Axelsen JB, Yaari R, Grenfell BT, Stone L (2014) Multiannual forecasting of seasonal influenza

449 dynamics reveals climatic and evolutionary drivers. *PNAS*, **111**, 9538–9542.

450 Breiman L (1999) Random Forests. *Machine Learning*, **45**, 1–35.

451 Bürgi HR, Bührer H, Keller B (2003) Long-Term Changes in Functional Properties and

452 Biodiversity of Plankton in Lake Greifensee (Switzerland) in Response to Phosphorus

453 Reduction. *Aquatic Ecosystem Health & Management*, **6**, 147–158.

454 Carey CC, Ibelings BW, Hoffmann EP, Hamilton DP, Brookes JD (2012) Eco-physiological

455 adaptations that favour freshwater cyanobacteria in a changing climate. *Water Research*, **46**,

456 1394–1407.

457 Chorus I, Bartram J (1999) *Toxic Cyanobacteria in Water: A Guide to Their Public Health*

458 *Consequences, Monitoring, and Management*. E & FN Spon, 416 pp.

459 Conley DJ, Paerl HW, Howarth RW et al. (2009) Controlling eutrophication: Nitrogen and

460 Phosphorus. *Science*, **323**, 1014–1015.

461 Demers S, Roy S, Gagnon R, Vignault C (1991) Rapid light-induced changes in cell

462 fluorescence and in xanthophyll-cycle pigments of Alexandrium excavatum (Dinophyceae) and

463 Thalassiosira pseudonana (Bacillario-phyceae): a photo-protection mechanism . *Marine Ecology*

464 *Progress Series*, **76**, 185–193.

465 Edwards KF, Thomas MK, Klausmeier CA, Litchman E (2015) Light and growth in marine

466 phytoplankton: allometric, taxonomic, and environmental variation. *Limnology and*

467 *Oceanography*, **60**, 540–552.

468    Edwards KF, Thomas MK, Klausmeier CA, Litchman E (2016) Phytoplankton growth and the

469    interaction of light and temperature: A synthesis at the species and community level. *Limnology*

470    *and Oceanography*, **61**, 1232–1244.

471    Eilers PHC, Peeters JCH (1988) A model for the relationship between light intensity and the rate

472    of photosynthesis in phytoplankton. *Ecological Modelling*, **42**, 199–215.

473    Ewers RM, Andrade A, Laurance SG, Camargo JL, Lovejoy TE, Laurance WF (2017) Predicted

474    trajectories of tree community change in Amazonian rainforest fragments. *Ecography*, **40**, 26–

475    35.

476    Falkowski PG, Barber RT, Smetacek V (1998) Biogeochemical controls and feedbacks on ocean

477    primary production. *Science*, **281**, 200–206.

478    Field CB, Behrenfeld MJ, Randerson JT, Falkowski PG (1998) Primary production of the

479    biosphere: Integrating terrestrial and oceanic components. *Science*, **281**, 237–240.

480    Goldman JC, Glibert PM (1982) Comparative rapid ammonium uptake by four species of marine

481    phytoplankton. *Limnology and Oceanography*, **27**, 814–827.

482    Hemme D, Veyel D, Mühlhaus T et al. (2014) Systems-Wide Analysis of Acclimation Responses

483    to Long-Term Heat Stress and Recovery in the Photosynthetic Model Organism

484    Chlamydomonas reinhardtii. *The Plant Cell*, **26**, 4270–4297.

485    Houlahan JE, McKinney ST, Anderson TM, McGill BJ (2017) The priority of prediction in

486    ecological understanding. *Oikos*, **126**, 1–7.

487    Hunter-Cevera KR, Neubert MG, Solow AR, Olson RJ, Shalapyonok A, Sosik HM (2014) Diel

488    size distributions reveal seasonal growth dynamics of a coastal phytoplankter. *PNAS*, **111**,

489    9852–7.

490    Hunter-Cevera KR, Neubert MG, Olson RJ, Solow AR, Shalapyonok A, Sosik HM (2016)

491    Physiological and ecological drivers of early spring blooms of a coastal phytoplankter. *Science*,

492    **354**, 326–329.

493    Ishwaran H. and Kogalur U.B. (2017). Random Forests for Survival, Regression and

494    Classification (RF-SRC), R package version 2.4.2.

495    Ishwaran H. and Kogalur U.B. (2007). Random survival forests for R. R News **7**(2), 25-31.

496    Jochimsen MC, Kümmerlin R, Straile D (2013) Compensatory dynamics and the stability of

497    phytoplankton biomass during four decades of eutrophication and oligotrophication. *Ecology*

498    *Letters*, **16**, 81–89.

499    Kehoe M, O'Brien K, Grinham A, Rissik D, Ahern KS, Maxwell P (2012) Random forest algorithm

500    yields accurate quantitative prediction models of benthic light at intertidal sites affected by toxic

501    Lyngbya majuscula blooms. *Harmful Algae*, **19**, 46–52.

502    Kehoe MJ, Chun KP, Baulch HM (2015) Who Smells? Forecasting Taste and Odor in a Drinking

503    Water Reservoir. *Environmental Science and Technology*, **49**, 10984–10992.

504    Kremer CT, Williams AK, Finiguerra M et al. (2016) Realizing the potential of trait-based aquatic

505    ecology: New tools and collaborative approaches. *Limnology and Oceanography*, **62**, 253–271.

506    Litchman E, Klausmeier CA (2008) Trait-based community ecology of phytoplankton. *Annual*

507    *Review of Ecology, Evolution, and Systematics*, **39**, 615–639.

508    Magurran AE, Baillie SR, Buckland ST et al. (2010) Long-term datasets in biodiversity research

509    and monitoring: assessing change in ecological communities through time. *Trends in Ecology &*

510    *Evolution*, **25**, 574–582.

511    McGill BJ, Enquist BJ, Weiher E, Westoby M (2006) Rebuilding community ecology from

512    functional traits. *Trends in Ecology & Evolution*, **21**, 178–85.

513    Mouquet N, Lagadeuc Y, Devictor V et al. (2015) Predictive ecology in a changing world. *Journal*

514    *of Applied Ecology*, **52**, 1293–1310.

515    Müller R, Pfeifroth U, Träger-Chatterjee C, Cremer R, Trentmann J, Hollmann R (2015) Surface

516    Solar Radiation Data Set - Heliosat (SARAH) - Edition 1.

517    Paerl HW, Hall NS, Calandrino ES (2011) Controlling harmful cyanobacterial blooms in a world

518    experiencing anthropogenic and climatic-induced change. *Science of the Total Environment*,

519    **409**, 1739–1745.

520    Paerl HW, Huisman J (2009) Climate change: A catalyst for global expansion of harmful

521    cyanobacterial blooms. *Environmental Microbiology Reports*, **1**, 27–37.

522    Petchey OL, Pontarp M, Massie TM et al. (2015) The ecological forecast horizon, and examples

523    of its uses and determinants. *Ecology Letters*, **18**, 597–611.

524    Pomati F, Jokela J, Simona M, Veronesi M, Ibelings BW (2011) An automated platform for

525    phytoplankton ecology and aquatic ecosystem monitoring. *Environmental Science &*

526    *Technology*, **45**, 9658–65.

527    Pomati F, Matthews B, Jokela J, Schildknecht A, Ibelings BW (2012) Effects of re-

528    oligotrophication and climate warming on plankton richness and community stability in a deep

529    mesotrophic lake. *Oikos*, **121**, 1317–1327.

530    Pomati F, Kraft NJB, Posch T, Eugster B, Jokela J, Ibelings BW (2013) Individual cell based

531    traits obtained by scanning flow-cytometry show selection by biotic and abiotic environmental

532    factors during a phytoplankton spring bloom. *PLoS one*, **8**, e71677.

533    Rice EW, Baird RB, Eaton AD, Clesceri LS (eds.) (2012) *Standard Methods for the Examination*

534    *of Water and Wastewater*, 22nd edn. American Water Works Association/American Public

535    Works Association/Water Environment Federation.

536    Rivero-Calle S, Gnanadesikan A, Del Castillo CE, Balch WM, Guikema SD (2015) Multidecadal

537    increase in North Atlantic coccolithophores and the potential role of rising CO2. *Science*, **350**,

538    1533–1537.

539    Schulz J, Albert P, Behr H-D et al. (2008) Operational climate monitoring from space: the

540    EUMETSAT satellite application facility on climate monitoring (CM-SAF). *Atmospheric Chemistry*

541    *and Physics Discussions*, **8**, 8517–8563.

542    Schwaderer AS, Yoshiyama K, de Tezanos Pinto P, Swenson NG, Klausmeier CA, Litchman E

543    (2011) Eco-evolutionary differences in light utilization traits and distributions of freshwater

544    phytoplankton. *Limnology and Oceanography*, **56**, 589–598.

545    Smith VH (1983) Low nitrogen to phosphorus ratios favor dominance by blue-green algae in lake

546    phytoplankton. *Science*, **221**, 669–671.

547     R Core Team (2017). R: A language and environment for statistical computing. R Foundation

548    for Statistical Computing, Vienna, Austria. https://www.R-project.org/.

549    Thomas MK, Kremer CT, Klausmeier CA, Litchman E (2012) A global pattern of thermal

550    adaptation in marine phytoplankton. *Science,* **338**, 1085–1088.

551    Thomas MK, Kremer CT, Litchman E (2016) Environment and evolutionary history determine the

552    global biogeography of phytoplankton temperature traits. *Global Ecology and Biogeography*, **25**,

553    75–86.

554    Thomas MK, Aranguren-Gassis M, Kremer CT, Gould MR, Anderson K, Klausmeier CA,

555    Litchman E (2017) Temperature-nutrient interactions exacerbate sensitivity to warming in

556    phytoplankton. *Global Change Biology*, **23**, 3269–3280.

557    Williams JW, Jackson ST, Kutzbach JE (2007) Projected distributions of novel and disappearing

558    climates by 2100 AD. *PNAS*, **104**, 5738–42.

559    Zhu K, Chiariello NR, Tobeck T, Fukami T, Field CB (2016) Nonlinear, interacting responses to

560    climate limit grassland production under global change. *PNAS*, **113**, 10589–10594.

561    Zimmer A, Katzir I, Dekel E, Mayo AE, Alon U (2016) Prediction of multidimensional drug dose

562    responses based on measurements of drug pairs. *PNAS*, **113**, 10442–10447.

563

564

565

**Fig. 1. Dynamics of cell density of the total phytoplankton community, in both the high-frequency and long-term datasets from Greifensee. High-frequency measurements were made every 4 hours in summer-fall 2014 and 2015, at six depths. Long-term measurements were made monthly from 1984 to 2016 and were integrated over the t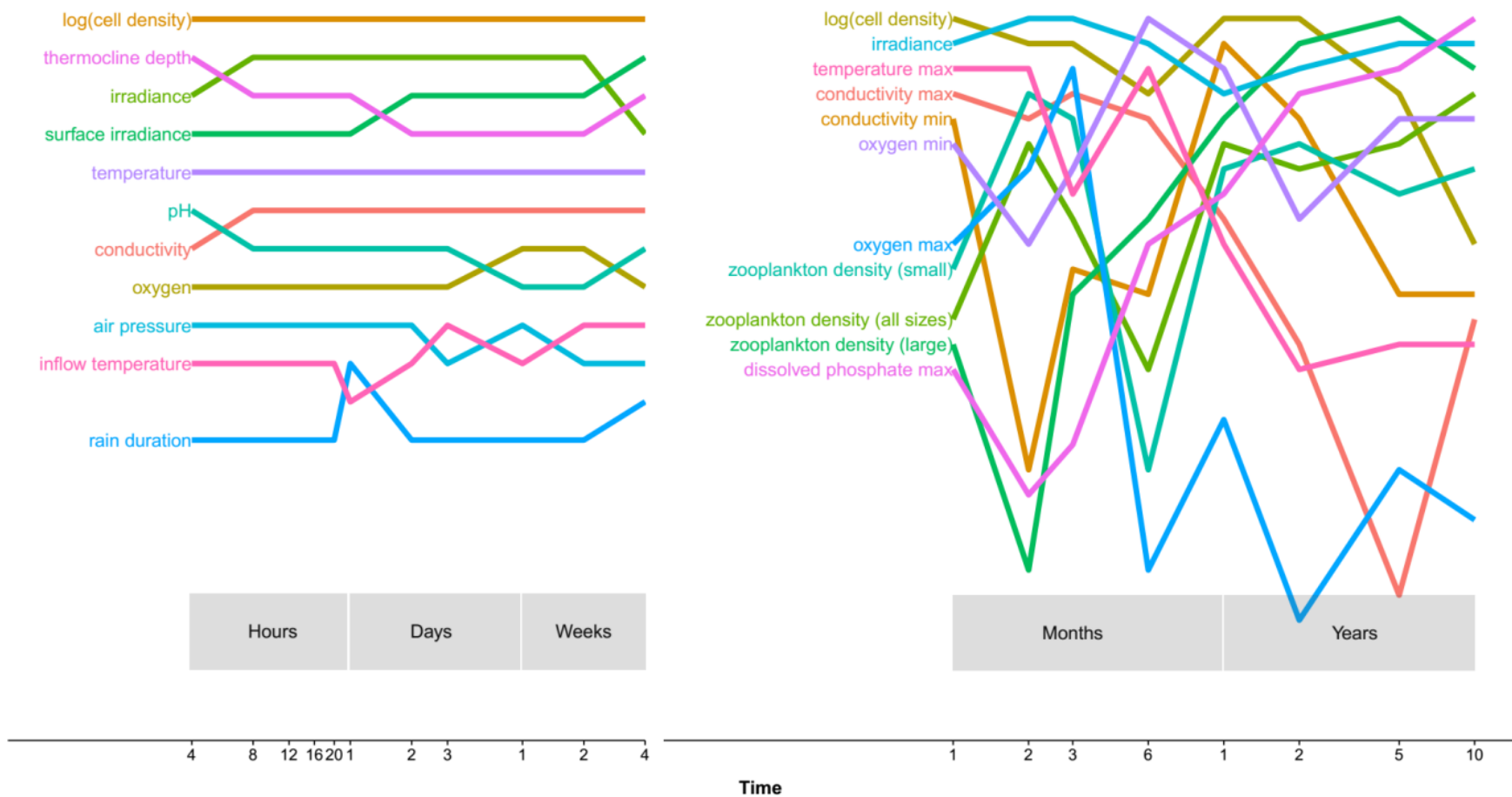op 20m. Note that X-axes are on different scales in each panel. Y-axes are identical for the top two panels but differ for the third.**

Fig. 2. Dynamics of cell density of the cyanobacteria and eukaryotic phytoplankton, in both the high-frequency and long-term datasets from Greifensee. High-frequency measurements were made every 4 hours in summer-fall 2014 and 2015, at six depths. Long-term measurements were made monthly from 1984 to 2016 and were integrated over the top 20m. Note that X- and Y-axes are on different scales in each column.

**Fig. 3. Decline in predictability of the phytoplankton community with time, characterised by the random forest pseudo-$R^2$. The predictive contribution of environmental information increased with increasing time lag (distance between solid and dashed lines increases). Cyanobacteria were consistently more predictable than eukaryotes. Despite overlap between high-frequency and long-term datasets at a time lag of 1 month, there is a decline in predictability likely driven by the lack of depth resolution in plankton and environmental data in the long-term dataset. The spike in $R^2$ of the 'cell density only' models at 1 year reflects strong annual cycles in density. Note that pseudo-$R^2$ values can go negative (see Methods), and we rounded a few slightly negative values up to zero. We present the same results in terms of change in Mean Absolute Error with increasing time lag in Fig. S5.**

**Fig. 4. The most important predictors of phytoplankton cell density at different time lags, ordered by descending rank. Light and temperature (directly, or indirectly through thermocline depth) were important predictors at most time scales. In the long-term dataset, phosphorus and zooplankton density become highly important predictors at time scales of >1 year. Only the most important variables are shown here, for legibility (the top 5 predictors contribute >80% of the predictive power in most cases). In the high-frequency dataset, only variables that are in the top 10 most important for at least one time lag are shown, while in the long-term dataset we show only variables that appear in the top 5 most important at least once. See Tables S4 and S5 for the importance of all variables tested.**

600