

# 1 Adaptive pseudo-real-time forecasting of phytoplankton communities

2 Trevor Page<sup>1\*\*</sup>, Paul J Smith<sup>1,3</sup>, Keith J Beven<sup>1</sup>, Ian D Jones<sup>2</sup>, J Alex Elliott<sup>2</sup>, Stephen C Maberly<sup>2</sup>, Eleanor B  
3 Mackay<sup>2</sup>, Mitzi De Ville<sup>2</sup> and Heidrun Feuchtmayr<sup>2</sup>.

4 <sup>1</sup> Lancaster Environment Centre, Library Avenue, Lancaster University, Lancaster, LA1 4YQ, UK.

5 <sup>2</sup> Lake Ecosystems Group, Centre for Ecology & Hydrology, Lancaster Environment Centre, Library  
6 Avenue, Bailrigg, Lancaster, LA1 4AP, UK.

7 <sup>3</sup> ECMWF, Shinfield Park, Reading, RG2 9AX, UK

8 \*\* Corresponding Author: [t.page@lancaster.ac.uk](mailto:t.page@lancaster.ac.uk)

9 **Key words:** Phytoplankton model, forecasting, data assimilation, Ensemble Kalman Filter,  
10 cyanobacteria, PROTECH.

## 11 Abstract

12 Evaluation of the potential for forecasting of algal blooms using the phytoplankton community model  
13 PROTECH was undertaken in pseudo-real-time. This was achieved within a data assimilation scheme using  
14 the Ensemble Kalman Filter to allow uncertainties and model nonlinearities to be propagated to forecast  
15 outputs. Testing was done on two mesotrophic lakes in the English Lake district, which have differing  
16 depths and nutrient regimes. Some forecasting success was shown for chlorophyll *a*, but not all forecasts  
17 were able to perform better than a persistence forecast. There was a general reduction in forecast skill  
18 with increasing forecasting period but forecasts for up to four or five days showed noticeably greater  
19 promise than those for longer periods. Associated forecasts of phytoplankton community structure were  
20 broadly consistent with observations but their translation to cyanobacteria forecasts is challenging owing  
21 to functional similarities between species which may or may not be cyanobacteria.

22 . It was concluded that higher frequency in-lake chlorophyll  $a$  and nutrient observations should help to  
23 improve forecasts but it remains to be seen how far the forecasting system can be used to identify algal  
24 bloom conditions in this type of lake.

## 25 **1 Introduction**

26 Algal blooms are a global problem affecting water resources, recreation and ecosystems (Carmichael,  
27 1992; Smith, 2003; World Health Organization, 1999). These problems are particularly acute when  
28 cyanobacterial species dominate because of the risk of toxin production that can cause adverse effects to  
29 humans and wildlife (Metcalf and Codd, 2009). In addition, water supply companies face associated  
30 problems such as poor taste and odour and, in extreme cases, high concentrations of algal-derived toxins  
31 which are costly to manage (Pretty *et al.*, 2003; Dodds *et al.*, 2009; Michalak, 2016). Costs associated with  
32 implementation of management strategies are growing because of increased bloom frequency (Ho and  
33 Michalak, 2015) because of the effects of nutrient enrichment and climate change (Paerl and Huisman,  
34 2008; Brookes and Carey, 2011; Rigosi *et al.* 2014). As a result, there is an urgent need for reliable  
35 predictions of algal bloom formation to enable timely management interventions.

36 Forecasting algal blooms in lakes is relatively new (Kim *et al.*, 2014) but is increasingly becoming a  
37 requirement for lake and reservoir managers (Huang *et al.*, 2013; Recknagel *et al.* 2014; Xiao *et al.*, 2017)  
38 to help inform decisions regarding the most cost-effective management strategies. The fact that  
39 limnology is rapidly becoming data-rich (Marcé *et al.*, 2016; Xiao *et al.*, 2014) means that effective real-  
40 time forecasts are increasingly more feasible. However, forecast simulations will be inherently uncertain  
41 for a number of reasons including input data resolution and simplifications in model process  
42 representation. These uncertainties will have implications for the accuracy and reliability of a forecast and  
43 therefore effort is required to allow for modelling uncertainty. Data assimilation (DA) is one approach to  
44 reducing forecast uncertainty, but has, to date, received relatively little attention for forecasting

45 phytoplankton community dynamics. There is hence a need to test different DA methodologies across  
46 different lake systems and different models.

47 There are still relatively few studies for operational lake forecasting systems and various approaches have  
48 been taken such as using: Ensemble Kalman Filter (EnKF; Evensen, 1994) schemes and physically-based  
49 simulation models (e.g. Allen *et al.*, 2003, Huang *et al.* 2013 and Kim *et al.*, 2014); evolutionary  
50 computation (Recknagel *et al.*, 2014; Ye *et al.*, 2014); Lagrangian particle tracking model methods (Rowe  
51 *et al.*, 2016); and using a combination of wavelet analysis and neural networks (Luo *et al.*, 2011; Xiao *et*  
52 *al.*, 2017). The EnKF has been developed to deal with highly non-linear model dynamics which cannot be  
53 represented well using the traditional Kalman Filter. Phytoplankton population dynamics are highly non-  
54 linear with multiple modes of behaviour that can respond rapidly to threshold-type effects and are prone  
55 to rapid changes in their physical and chemical environment (e.g. water temperature, light levels and  
56 available nutrients). This makes the EnKF a suitable choice to exploring algal bloom forecasting when  
57 coupled with a phytoplankton community model.

58 Here we assess our ability to make pseudo-real-time forecasts of phytoplankton communities in two lakes  
59 in the north west of England, which are prone to cyanobacteria blooms during the summer. Forecasts are  
60 made using a modified version of the phytoplankton community model PROTECH (Reynolds *et al.*, 2001)  
61 within a DA scheme using the EnKF. The version of PROTECH employed is appropriate for this problem as  
62 it is intermediate in its complexity between physically-based coupled 3-dimensional hydrodynamic-  
63 biochemical models and more simplistic “black box models” which have both been used in this context.  
64 More complex models are extremely computationally expensive in forecasting (Huang *et al.*, 2012;  
65 Recknagel, *et al.*, 2014), such that only a limited number of ensemble members can be used (Kim *et al.*,  
66 2014); and simple black box models may not be able to represent phytoplankton community dynamics

67 driven by ecological strategies that are represented in phytoplankton community models such as  
68 PROTECH.

69 We aim to determine the efficacy of phytoplankton community forecast simulations, evaluate the EnKF  
70 as a DA strategy and investigate the ensemble size required for making consistent forecasts. Ultimately,  
71 success will rely on the modelling strategy being sufficiently effective to capture the necessary short-term  
72 phytoplankton community dynamics, given the available meteorological forecasts and limitations  
73 associated with driving data. Demonstrating the efficacy of the approach therefore requires a robust  
74 appraisal procedure with predictions tested qualitatively and quantitatively against appropriate  
75 benchmarks. This approach allows other pertinent questions to be investigated; namely, how does  
76 forecasting reliability diminish with time-scale of forecast and, most pertinently, what can be learnt from  
77 any forecasting failure regarding future model development and optimisation of monitoring strategies.

## 78 **2 Methods**

### 79 **2.1 Study lakes**

80 This study considers two lakes in the English Lake District of North West England with differing depths and  
81 nutrient regimes (Table 1). The catchments associated with each of the lakes are predominantly hill land,  
82 rough-grazed by sheep throughout the year and contain towns and villages that are tourist destinations  
83 and are hence associated with seasonal increases in lake nutrient inputs. Windermere is England's largest  
84 lake and comprises two basins connected at a shallow region approximately halfway along its main axis.  
85 The two basins are usually considered separately as they have different characteristics: both basins are  
86 monomictic and mesotrophic; the south basin was modelled in this study. Esthwaite Water is a small,  
87 generally monomictic and occasionally dimictic, lake that has been subject to eutrophication for many  
88 decades because of elevated phosphorus levels (Bennion *et al.*, 2000; Dong *et al.*, 2012): cyanobacterial  
89 blooms are common in the summer to early autumn. Previous work has found that internal sources form

90 an important component of the P budget of the lake (Hall *et al.* 2000; Heaney *et al.*, 1992 and Mackay *et*  
91 *al.*, 2014).

92 **Table 1 Study Lakes and primary characteristics<sup>§</sup>**

Name/location	Mean Depth (m)	Max. Depth (m)	Max. Length (m)	Volume (m <sup>3</sup> )	Catchment Area (km <sup>2</sup> )	Residence Time (days)
Windermere (South Basin)	16.8	41	9300	1.06 x 10 <sup>8</sup>	230.5	100
Esthwaite Water	6.4	15.5	2500	5.97 x 10 <sup>6</sup>	17.1	100

93 <sup>§</sup> Details from Ramsbottom (1976)

## 94 2.2 Data

### 95 2.2.1 Forcing inputs: meteorological forecasts

96 The primary forcing inputs were meteorological forecasts provided by the European Centre for Medium-  
97 term Weather Forecasts (ECMWF) Ensemble Prediction System. The 10-day-ahead forecasts include an  
98 ensemble of 50 simulations from perturbed initial states (at 32 km<sup>2</sup> resolution) and stochastic  
99 perturbations of model parameters (see Buizza *et al.*, 1999 and Ollinaho *et al.*, 2016). The re-initialisation  
100 of model states in the ECMWF forecasting system is implemented using a higher resolution 3-hour  
101 forecast each day. As this re-initialisation is repeated each day, and as perturbations are random, there is  
102 no specific relationship between individual ensemble members in subsequent days. The forecast  
103 associated with each ensemble member was hence treated as independent from prior forecasts for this  
104 study. Daily averages of forecasts were used (i.e. the average of 3-hourly forecasts for days 1-6 and of 6-  
105 hourly forecasts day 6-10) for consistency with the daily timestep of PROTECH. Historic forecasts were  
106 obtained for 2008, 2009 and 2010 and used in pseudo-real-time. Given the scale of the forecast grid, each  
107 forecast variable was “downscaled” to local data as described in the next section.

108

### 109 2.2.2 Sampling meteorological forecasts

110 Downscaling relationships were developed for air temperature, wind speed, precipitation, cloud cover,  
111 relative humidity and solar radiation (Table 2). For air temperature a relationship was identified between  
112 forecasted temperatures and observed temperatures using linear regression. Residuals from this initial  
113 analysis helped identify an additional hysteretic relationship between forecasted and observed  
114 temperatures, which was attributed to a lake thermal effect; this effect was implemented as an additional  
115 correction for each day of the year. Similarly, wind speed was corrected using a linear correction factor  
116 coupled with an additional correction based upon wind direction; this was required owing to complex  
117 mountainous topography and lake-axis orientation. A wind-rose with sectors of 30 degrees was used to  
118 classify forecasted wind speeds and a sector-specific correction was applied. The uncertainty associated  
119 with the corrections was represented by fitting a gamma distribution to the data in each sector. All other  
120 variables (precipitation, cloud cover, relative humidity and solar radiation), were corrected using a  
121 correction multiplier identified using linear regression, without propagating the uncertainty in the  
122 relationship. The uncertain relationships for air temperature and wind speed were resampled as  
123 perturbations of the ensemble members allowing investigation of the effect of different ensemble sizes.

### 124 **2.2.3 Nutrient Inputs**

125 Knowledge of diffuse nutrient inputs for the study lakes is relatively poor. Observations available were  
126 from approximately monthly frequency routine monitoring and did not cover all river inputs. Both lakes  
127 are also impacted by point sources from waste water treatment works (WwTW) and Esthwaite is subject  
128 to significant internal P fluxes (Mackay *et al.*, 2014). Diffuse nutrient inputs and WwTW inputs (where  
129 included) were treated as reported by Page *et al.* (2017) and these inputs were modified by a  
130 multiplicative parameter included in the EnKF scheme (Table 3). For Windermere, upstream lake inputs  
131 of nutrients (and chlorophyll *a*) were treated as reported by Page *et al.* (2017) but were not included in  
132 the EnKF scheme.

133 **Table 2 Forcing inputs and downscaling relationships**

Model Inputs	Downscaling factor/relationship	Uncertainty sampled
Air Temp ( $T_a$ ; K)	Windermere: $0.095(T_a^s) + 279.75^{**}$ Esthwaite Water: $0.013(T_a^s) + 280.16^{**}$	Y (Regression)
Solar Radiation (SR; $Wm^{-2}$ )	0.85	N
Wind Speed (W; $m s^{-1}$ )	$0.38^{\$}$	Y (Gamma Dist.)
Relative Humidity (RH; %)	1	N
Cloud Cover (Cc; eighths)	1.25	N
Rainfall (R; mm)	3	N
Nutrient Inputs (P; N; $SiO_2/ mg m^{-3}$ )	Section 2.2.3	Y (Gamma Dist.)

134  $T_a^s$  is the forecast air temperature (K); \*\* see Section 2.2.2 for additional lake-effect correction; \$ see Section 2.2.2 for additional  
135 wind direction correction.

136 **2.2.4 Data for assimilation**

137 Specific years where the observed data were of the highest frequency, were chosen to test the DA  
138 strategy. High frequency data from the automatic lake monitoring systems (Madgwick *et al.*, 2006;  
139 Mackay *et al.*, 2014) were available and were aggregated to daily values. The variables used for DA are  
140 listed in Table 3. The “observed” temperatures for the epilimnion ( $T_e$ ) and hypolimnion ( $T_h$ ) used to  
141 compare with the modelled variables for these layers were calculated as volume-weighted averages of  
142 thermistor chain data, using the simulated epilimnetic depth to delineate the hypolimnion and epilimnion.  
143 The “observed” epilimnetic depth ( $D_e$ ) was estimated using a density gradient method (e.g. see Read *et*  
144 *al.*, 2011). In addition to the automatic monitoring, routine monitoring was carried out at the buoy  
145 location at a frequency of approximately every 14 days and included chlorophyll *a*, soluble reactive  
146 phosphorus (SRP), dissolved inorganic nitrogen (DIN) and silica ( $SiO_2$ ) (Table 3). These observations were  
147 derived from a water sample at the buoy location integrated over 0-7 m depth (Windermere) or 0-5 m  
148 depth (Esthwaite Water) (Maberly *et al.*, 2010).

149

150 **2.3 Modelling methodology**

151 The modelling strategy employed was designed to represent the different facets of the forecasting  
152 system as simply as possible to reduce computational burden, whilst retaining the requirement to

153 explicitly simulate phytoplankton community structure and, specifically, to estimate the likely  
 154 concentrations of cyanobacteria from this structure. Thus, the catchment-lake system was simulated  
 155 using a suite of models of differing complexity from purely data-based (statistically estimated) transfer  
 156 function (TF) models and processed-based models which are consistent, in their complexity, with the  
 157 available data. A schematic of how the models were combined in the forecasting system is presented in  
 158 Figure 1 and each model is described in this section. The modelling system is structured around the  
 159 rationale that epilimnetic depth must be estimated as accurately as possible so that the phytoplankton  
 160 model, PROTECH, is more likely to provide good estimates of phytoplankton community structure; in  
 161 PROTECH, community structure is simulated using functional algal types as classified by Reynolds (1988)  
 162 as outlined in the next section. The simple conceptual model that estimates epilimnetic depth is a heat  
 163 energy “balance” model that requires estimates of epilimnetic temperature and energy fluxes to the  
 164 epilimnion, including those associated with river inflows and outflows.

165 **Table 3 Observed data assimilated in the EnKF scheme**

Assimilated state	Frequency	Source
Epilimnetic Temperature (°C)	Daily	buoy obs.
Hypolimnetic Temperature (°C)	Daily	buoy obs.
Epilimnetic depth (m)	Daily	buoy obs.
Chlorophyll a (mg m <sup>-3</sup> )	≈14 days	Monitoring
Nutrient Inputs (SRP; N; SiO <sub>2</sub> / mg m <sup>-3</sup> )	≈14 days	Monitoring

166

167 The TF models, epilimnetic depth model and PROTECH are run sequentially; the TF and epilimnetic depth  
 168 models provide forecast estimates of river flow, epilimnetic depth, epilimnetic temperature and  
 169 hypolimnetic temperature as inputs to PROTECH. Data assimilation is employed for the two primary  
 170 models (the epilimnetic depth model and PROTECH) using two separate EnKF schemes that assimilate  
 171 observations at different intervals; the epilimnetic depth model scheme assimilates epilimnetic depth and  
 172 epilimnetic temperature estimates as well as hypolimnetic temperature estimates on a daily basis and the  
 173 scheme for PROTECH assimilates nutrient and chlorophyll *a* concentrations approximately every 14 days.



174 **2.3.1 The PROTECH model**

175 PROTECH (Reynolds *et al.*, 2001) is a lake phytoplankton community model that runs on a daily time-step.  
176 It is a 1-dimensional model where the lake is represented by horizontal layers. In the model representation  
177 all layers are assumed to be fully mixed throughout the epilimnion. River inputs drive fluxes of diffuse  
178 nutrients as well as the flushing of phytoplankton. Upstream lake inputs are treated as river inputs but  
179 are given the phytoplankton concentrations associated with the upstream lake, where data are available.

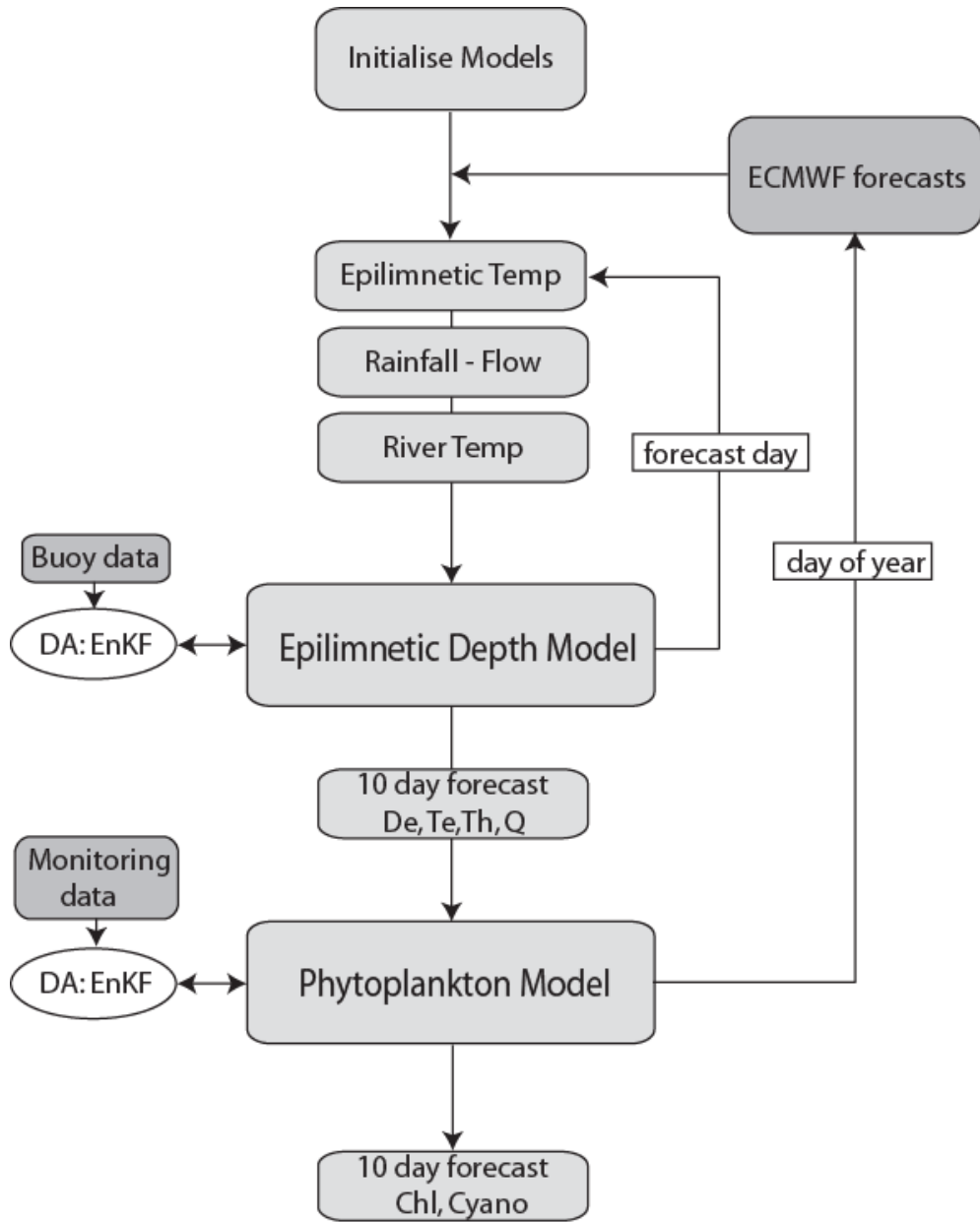
180 Underwater light for model layer  $i$  is calculated using:

181 
$$l_i = I_{surf} \cdot e^{(-\varepsilon \cdot d_i)} \quad (1)$$

182

183 Where:  $I_{surf}$  is the daily surface light flux,  $d$  is the depth from the lake surface,  $\varepsilon$  is the light extinction  
184 coefficient resulting from the sum of lake-specific abiotic water attenuation ( $\varepsilon_b$ ) and the extinction of light  
185 associated with the concentration of phytoplankton at each timestep multiplied by the parameter  $\varepsilon_a$ . In  
186 the layers from the surface to the epilimnetic depth, the available light is represented by the geometric  
187 mean of the epilimnetic layers and hence assumes that phytoplankton spend an equal time in each layer  
188 at each timestep. Phytoplankton population dynamics are simulated using the following equation which  
189 describes the change in chlorophyll  $a$  concentration ( $X$ ) of each phytoplankton species selected to  
190 represent the algal community (Reynolds 2001):

191 
$$\frac{\Delta X}{\Delta t} = (r' - S - G - F) \cdot X \quad (2)$$



192

193

194

195

196

*Figure 1. Schematic diagram of the forecasting system. The schematic shows sequential model input-output structure and DA strategy. De is epilimnetic depth; Te is epilimnetic temperature; Th is hypolimnetic temperature, Q is lake inflow/outflow and Chl and Cyano are the concentration of total phytoplankton chlorophyll a and cyanobacterial chlorophyll a respectively.*

197 where  $r'$  is the growth rate,  $S$  is settling loss,  $G$  is a grazing loss and  $F$  is the loss due to flushing. The growth  
198 rate is defined for each layer using:

$$199 \quad r' = \min \{r'_{(\theta)}, r'_{(P)}, r'_{(N)}, r'_{(Si)}\} \quad (3)$$

200 where  $r'_{(\theta,l)}$  is the growth rate at a given temperature ( $\theta$ ) and daily photoperiod ( $l$ ) and  $r'_P, r'_N, r'_{Si}$  are the  
201 growth rates determined by phosphorus, nitrogen and silica concentrations. The final growth rate ( $r'_{corr(\theta,l)}$ )  
202 is a corrected rate allowing for dark respiration using equation 4. This is required as the model growth  
203 equations are net of basal metabolism but not dark respiration burden.

$$204 \quad r'_{corr(\theta,l)} = R_{d(\theta)} \cdot r'_{(\theta,l)} - (1 - R_{d(\theta)}) \cdot r'_{(\theta,l)} \quad (4)$$

205 Where  $R_{d(\theta)}$  is the dark respiration rate at temperature  $\theta$ . The phytoplankton used for this study are  
206 presented in Table Supp. 2.

207 PROTECH simulates the dynamics of the species chosen to represent the algal community of a given lake.  
208 Species are represented by their morphology, nutrient requirements (i.e. silica requirement and nitrogen  
209 fixing ability) and their vertical movement strategies. The number of species simulated is nominally eight  
210 (although unlimited) and they are chosen to represent the dominant functional types of the system.  
211 Simulations hence represent the behaviour of the functional algal community rather than the dynamics  
212 of specific species. The C-S-R functional phytoplankton classification of Reynolds (1988) is used to classify  
213 phytoplankton into morphologically defined groups relating to broad ecological strategies. The primary  
214 groups are: C-types, which are invasive, ecological pioneers that are small with high surface-to-volume  
215 ratios (e.g. *Chlorella*, and *Plagioselmis*); S-types which are 'stress tolerators' that tolerate relatively low  
216 nutrient availability and strong stratification (e.g. *Woronichinia*, *Microcystis* and *Oocystis*); and R-types  
217 which can harvest sufficient light at low levels to be able to maintain growth and are hence tolerant of  
218 well-mixed, intermittently insolated environments (e.g. *Asterionella*, *Aulacoseira* and *Oscillatoria*). Also

219 important for the lakes studied here, are CS-types, whose characteristics are intermediate between those  
 220 of C and S species (e.g. Dolichospermum, Aphanizomenon and Ceratium) and CSR-types (e.g.  
 221 Cryptomonas) that are intermediate between C-, S- and R-types.

### 222 2.3.2 Epilimnetic depth model

223 As a way of reducing computational burden, a simplified representation of lake thermal structure was  
 224 employed to estimate epilimnetic depth ( $D_e$ ). The simplified model works on the basis of *independent*  
 225 estimates of epilimnetic temperature and lake heat energy fluxes. The estimate of epilimnetic  
 226 temperature ( $T_e$ ) uses a TF model (see Section 2.3.3) with inputs of air temperature ( $T_a$ ), solar radiation,  
 227 wind speed ( $W_s$ ) and  $D_e$ . Air temperature solar radiation and wind speed are derived from the forecasts  
 228 and  $D_e$  estimates are from the previous simulation timestep. The independent estimates of heat energy  
 229 fluxes are calculated using the PROTECH energy flux function (see Reynolds *et al.*, 2001) for each timestep  
 230 using  $T_e$ , river temperature and flow magnitude, day length, cloud cover,  $T_a$ , Relative Humidity and  $W_s$ .

231 These two independent estimates are “balanced” to obtain hypolimnetic volume ( $V_h$ ) using:

$$232 \quad V_h = \frac{E_{\Delta T}}{\Delta T \cdot C_w \cdot \rho_w} \quad (5)$$

233 where,  $E_{\Delta T}$  is the heat energy associated with  $\Delta T$  (the difference between  $T_e$  and the hypolimnetic  
 234 temperature,  $T_h$ ),  $C_w$  is the specific heat capacity of water,  $\rho_w$  is the density of water. Equation 5 is solved  
 235 to find  $V_h$  where:  $\Delta T \cdot C_w \cdot \rho_w \cdot V_h \approx E_{\Delta T}$ . Subsequently, the epilimnetic volume ( $V_e$ ) and hence epilimnetic  
 236 depth ( $D_e$ ) are estimated by difference:

$$237 \quad V_e = V_t - V_h \quad (6)$$

238 where  $V_t$  is the total lake volume. The requirement for  $\Delta T$  is satisfied by calculating  $T_h$  using:

$$239 \quad T_h = \frac{E_{th}}{C_w \cdot \rho_w \cdot V_t} \quad (7)$$

240 where:  $E_{th}$  is the “background” heat energy in the lake (associated with  $T_h$  and  $V_t$ , as defined by Eqn. 7).  
241 During the forecast period,  $E_{th}$  remains at its previous value until updated during the data assimilation  
242 step. This treatment of  $E_{th}$  neglects the explicit downward transfer of energy from  $E_{\Delta T}$  to  $E_{th}$  for forecasting  
243 and assumes that these are negligible over this timescale: energy is, however, explicitly transferred  
244 downwards each time temperatures are updated during data assimilation. The sequence of calculations  
245 for each forecast timestep is:

- 246 1. Estimate lake surface temperature using TF model
- 247 2. Update  $E_{\Delta T}$ 
  - 248 I. Radiative energy fluxes
  - 249 II. River/upstream lake fluxes
    - 250 • Estimate river input volume using TF model
    - 251 • Estimate river temperature using TF model
    - 252 • Assume Upstream lake temperature = modelled lake temperature
  - 253 III. If  $E_{\Delta T} < 0$  lose energy from  $E_h$  (minimum energy set to 0°C)
- 254 3. Estimate  $T_h$  from  $E_{th}$
- 255 4. If  $E_{\Delta T} > 0$  and If  $T_e - T_h$  is greater than a threshold parameter (nominally set to 1°C) estimate  
256 epilimnetic depth by solving for the volume of water required to match  $E_{\Delta T}$  given  $\Delta T$ :  
257 subsequently estimate  $V_e$  and hence  $D_e$  by difference.

### 258 2.3.3 Transfer Function models

259 Transfer Function (TF) models were used for estimating lake surface temperature, river temperature and  
260 river inflows and outflows. Each model is a discrete-time TF identified directly from the available data.  
261 Both the model structures and parameters were identified using the Refined Instrumental Variable (RIV)  
262 algorithm (Young, 2015) implemented within the CAPTAIN Toolbox for Matlab™ (Taylor *et al.*, 2007). The

263 resulting model structures and parameter values are presented in Section (Supp. 1) and are either single  
264 input- or multi-input, single-output first order models of the general form:

$$265 \quad y_t = \frac{B_1(z-1)}{A(z-1)} U_1 + \frac{B_2(z-1)}{A(z-1)} U_2 + \dots + \frac{B_n(z-1)}{A(z-1)} U_n \quad (8)$$

266 where,  $y_t$  is the variable being estimated at time  $t$ ,  $U_{1-n}$  are model input vectors,  $A(z - 1)$  and  $B_n(z - 1)$   
267 are the model coefficients (polynomials in the backward shift operator: defined by  $y_t z^{-1} = y_{t-1}$ ) that  
268 number 1 to  $n$  in the case of  $B$  but note that in this form for MISO (multi-input single-output) TF the  
269 denominator ( $A$ ) is common to all  $n$  TF elements.

#### 270 **2.3.4 The Ensemble Kalman Filter**

271 The EnKF is a sequential Monte Carlo method which uses a stochastic ensemble of model simulations, and  
272 stochastic forcing, to propagate estimates of model states and (or) parameter values between assimilation  
273 timesteps. As the ensemble of model simulations is used in place of the linear propagation of an error  
274 covariance matrix (as in the traditional Kalman Filter), non-linear model dynamics are retained during  
275 model evolution and uncertainties are represented by the variation of the ensemble. When observations  
276 are available, each ensemble member is updated individually using a linear update equation (Eqn. 9) which  
277 relies on the assumption that the relationship between states and parameters can be described by  
278 multivariate Gaussian distributions. Rather than resampling the posterior distributions of the updated  
279 ensemble, the EnKF uses each updated ensemble member such that some of the non-Gaussian properties  
280 of the forecast are retained (Evenson, 2009). The procedure for the scheme is as follows:

281 1. The EnKF is initialised with an  $N$  number ensemble size, sampling states and parameters from *a priori*  
282 specified distributions (see below for specific details of this study) and  $N$  simulations for the forecast  
283 period are carried out. Where parameters are varied as part of the EnKF scheme, they are appended to  
284 the state matrix to give a state-parameter matrix.

285 2. When observed data are available for assimilation:

286 I. Apply a linear covariance inflation factor ( $I$ ) to each of the  $i$  states and parameters to reduce the  
287 tendency for low ensemble covariance and for spurious correlations associated with small  
288 ensemble size (Anderson, 2007; Anderson and Anderson, 1999; Evenson, 2009):

289  
290 
$$\varphi_{j,i}^a = I. (\varphi_{j,i}^a - \overline{\varphi_i^a}) + \overline{\varphi_i^a} \quad (9)$$

291

292 II. Generate  $N$  perturbations of the observations ( $Y$ ); it is essential that the uncertainty associated  
293 with the observations is sampled from a distribution with mean equal to the observed value and  
294 covariance ( $P^e$ ) to avoid bias in the update (Evenson, 2009) and to reduce further the tendency  
295 for the updated ensemble to have very low covariance (Moradkhani *et al.*, 2005).

296

297 III. Update the model states and parameters individually for the  $j^{th}$  ensemble member. This is done  
298 proportionally to the deviation of the states in the forecasted state-parameter matrix ( $\varphi^f$ ) from  
299 the vector of perturbed observations and the Kalman gain matrix ( $K$ ): note that the timestep  
300 suffix is omitted for clarity in the following equations:

301

302 
$$\varphi^a = \varphi^f + K(Y - H\varphi^f) \quad (10)$$

303 where,  $\varphi^a$  is the vector of updated states/parameters and  $H$  is a matrix that maps the model  
304 states to the observed sates. The appended parameters are updated using the cross-covariance  
305 between the predicted states and parameters. The Kalman gain matrix is calculated using:

306 
$$K = P_\varphi^f H^T (H(P_\varphi^f)H^T + P^e)^{-1} \quad (11)$$

307 where,  $P_{\phi}^f$  is the covariance matrix for the ensemble of forecasted state-parameter matrix.

308 IV. Apply any constraints on states and (or) parameter distributions (e.g. to keep them within  
309 physically reasonable ranges). This was implemented using a resampling scheme where if any  
310 state/parameter violated specified constraints (Table 4), the ensemble was resampled using a  
311 truncated distribution for that state/parameter in conjunction with a Gaussian copula to retain  
312 the ensemble's covariance structure.

313  
314 V. Make  $N$  number of simulations for the next forecast period using the updated state-parameter  
315 matrix.

#### 316 **2.3.5 Ensemble Kalman Filter scheme: Epilimnetic model**

317 As the epilimnetic model is very simple, all the main model states were used in the EnKF scheme. The  
318 states  $T_e$ ,  $T_h$  and  $D_e$  were updated using a daily assimilation frequency for the epilimnetic depth model.  
319 The “observed” values of these states are those estimated and described above.

#### 320 **2.3.6 Ensemble Kalman Filter scheme: PROTECH**

321 The choice of states and parameters included in the PROTECH EnKF scheme was made based on  
322 uncertainty and sensitivity analyses reported by Page *et al.* (2017). The Page *et al.*, study, which included  
323 the lakes studied here, identified that the main challenges for forecasting as uncertainties associated with:  
324 representing phytoplankton exposure to light and nutrient inputs (particularly phosphorus). The DA  
325 scheme was therefore defined to include the main model states, SRP, DIN, SiO<sub>2</sub> and chlorophyll  $a$ , as well  
326 the parameters associated with modifying nutrient inputs and underwater light (Table 4). These were  
327 updated at an approximately 14-day frequency set by the monitoring data. For Windermere both point  
328 source ( $WwTW_f$ ) and diffuse SRP inputs ( $P_{fact}$ ) parameters were included in the DA scheme; for Esthwaite



329 Water only the parameter modifying the diffuse SRP inputs was included as simulations which included a  
 330 simplified representation of sediment-derived SRP inputs did not provide improved results (these results  
 331 are not reported here).

332  
 333 To investigate the effect of ensemble size and to determine an acceptable ensemble size for the current  
 334 applications, ensemble member (EM) size was increased sequentially, using the scenarios EM50, EM100,  
 335 EM200, EM300 and EM400 (where the suffix is the size of the ensemble), until the forecast simulations  
 336 appeared consistent. These scenarios were generated by resampling the downscaled ECMWF forecast  
 337 distributions as described above and were used to force the suite of models used. For each of the forecast  
 338 scenarios, the error associated with the assimilated data and the variance inflation factors were  
 339 “optimised” manually to provide the best results. For consistency, and in the spirit of the pseudo-real time  
 340 treatment of the forecast simulations, the variance inflation factors were kept consistent across all lake-  
 341 years considered. For each of the assimilated variables, the variance was assumed to be proportional to  
 342 the magnitude of the variable of interest using a percentage. Additionally, a minimum variance was  
 343 applied to reduce the impact of very small observed values (e.g. where hypolimnetic SRP values are  
 344 observed to be very low or within the limit of detection) where the associated low variance would falsely  
 345 indicate low uncertainty.

346 **Table 4. States and parameters included in the ENKF scheme**

347

State/Parameter	Acceptable range	Observational error (%)	Initial distributions (uniform)**
Epilimnetic Temp. ( $T_e$ , °C)	2-25	5	5.5-7 (W); 4-6(E)
Hypolimnetic temp. ( $T_h$ , °C)	2-25	10	5.5-7 (W); 4-6(E)
Epilimnetic depth ( $D_e$ , m)	0.5-Lake depth	max. 5	41 (W); 15.5(E)
Chlorophyll <i>a</i> ( $\text{mg m}^{-3}$ )	$1e^{-6}$ - $1e^3$	10	3-4.5 (W); -4.5-6 (E)
Background light extinction ( $\epsilon_b$ , $m^{-1}$ )	0.15-0.9	N/A	0.15-0.6(W); 0.45-0.75(E)
Epilimnetic P conc. ( $P_e$ , $\text{mg m}^{-3}$ )	$1e^{-6}$ - $1e^4$	25	10-20(W); 8-15(E)
Epilimnetic DIN conc. ( $N_e$ , $\text{mg m}^{-3}$ )	$1e^{-6}$ - $1e^4$	25	400-700(W); 500-1100(E)
Epilimnetic SiO <sub>2</sub> conc. ( $Si_e$ , $\text{mg m}^{-3}$ )	$1e^{-6}$ - $1e^4$	25	1500-2500(W); 2000-2500(E)
Diffuse P input multiplier ( $P_f$ , dimensionless)	0.05-7	N/A	0.01-1.5

Diffuse DIN input multiplier ( $N_f$ , dimensionless)	0.1-3	N/A	0.5-1.2
Diffuse SiO <sub>2</sub> input multiplier ( $Si_f$ , dimensionless)	0.1-3	N/A	0.5-1.2
Point source P input multiplier ( $WwTW_f$ , dimensionless)	0.01-2	N/A	0.1-1.4

348 \*\* Where distributions are different for each lake W = Windermere; E = Esthwaite Water

### 349 2.3.7 Assessing forecast skill

350 Different studies have used different benchmarks to evaluate the goodness of fit of forecasts (*forecast*  
351 *skill*), which are often determined by their aims. Studies tend to use either some form of “reference”  
352 simulation or simulations that do not assimilate any observations (sometimes called “climatology”) which  
353 serve to quantify the DA effect (e.g. Allen *et al.*, 2003 and Kim *et al.*, 2014) or solely a measure of the  
354 goodness-of-fit to observations (e.g. the coefficient of determination,  $R_T^2$ ). Here, as our aim was to assess  
355 the value of the model for operational forecasting, we used a more stringent *persistence forecast* (e.g. see  
356 Stumpf *et al.*, 2009) which uses the most recent observations as the forecast for each *forecast timestep*  
357 until the next observation becomes available. In the sections below, the forecast skill was assessed using  
358 a persistence forecast for the entire annual timeseries and for the chlorophyll *a* forecast for which we  
359 have the most confidence in the observations. The goodness of fit of the benchmark and the simulated  
360 chlorophyll *a* forecasts are determined using the root-mean-square error (RMSE) as a measure. For the  
361 epilimnetic depth model, and other sub-models (i.e. TF models), goodness of fit is discussed more  
362 generally by comparison with observations using the coefficient of determination ( $R_T^2$ ). Assessment of  
363 the forecasts of phytoplankton community structure is made qualitatively as we have a significantly lower  
364 confidence in the absolute value of the observations.

## 365 3 Results and discussion

### 366 3.1 TF model results

367 Transfer function models were identified for epilimnetic temperature, river temperature and river inflows  
368 and outflows and all models provided good fits to the observed data during model identification:  $R_T^2$   
369 values were between 0.86 and 0.98 (Supp. Table 1). Model identification was carried out for the entire  
370 period of data available (see Supp. 1) such that they were not year specific models. As detailed above, in  
371 each case the models were used to forecast their respective variable deterministically.

## 372 **3.2 Forecasting epilimnetic depth and the phytoplankton community**

### 373 **3.2.1 Epilimnetic depth forecasts**

374 Epilimnetic depth forecast estimates were made for 2008-2010 for Windermere and 2008 and 2009 for  
375 Esthwaite Water within the parallel EnKF scheme. Although very simplistic, the epilimnetic depth model  
376 provided reasonable forecasts of epilimnetic depth when compared to those estimated from  
377 observations. For both lakes, the forecasts were stable and consistent using the smallest ensemble size of  
378 50 using a variance inflation factor of 1.25. Simulations for Windermere were better than for Esthwaite  
379 Water ( $R_T^2$  of 0.85 and 0.75 respectively for a 10-day-ahead forecast; Figs. 2a and 2b) and there were short  
380 periods with significant deviations from the 'observed' depths in both cases. Simulation of the timing of  
381 temporary stratification events at the beginning of the year was problematic for both lakes and  
382 simulations tended towards overly rapid mixing during autumn turnover, particularly for Esthwaite Water.  
383 Where significant deviations exist, they have the potential to reduce the forecast skill and therefore need  
384 to be improved, although, importantly, epilimnetic depth estimates for much of the high cyanobacterial  
385 bloom risk periods (i.e. during periods of strongest stratification) are reasonable. Given these results, the  
386 epilimnetic depth estimates for Windermere appear to be adequate out to 10-days-ahead but for  
387 Esthwaite they appear to be adequate for a much shorter lead time; for example, the 3-day-ahead forecast  
388 is a much better fit with an improved  $R_T^2$  of 0.81 (Fig. 2c). The adequacy of these estimates is assessed

389 more formally in association with the phytoplankton forecasts in comparison to the persistence forecast  
390 in the next section.

391

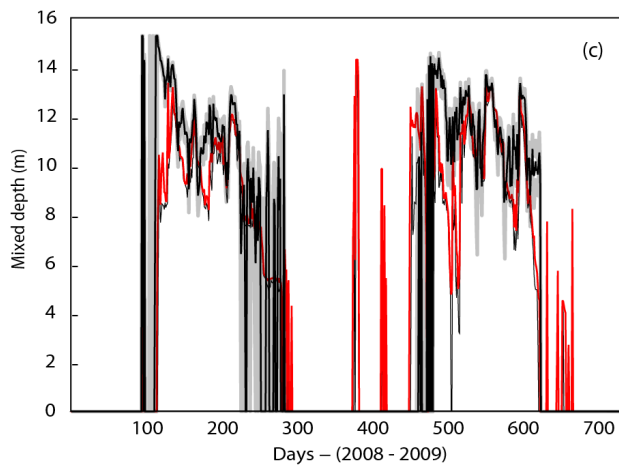
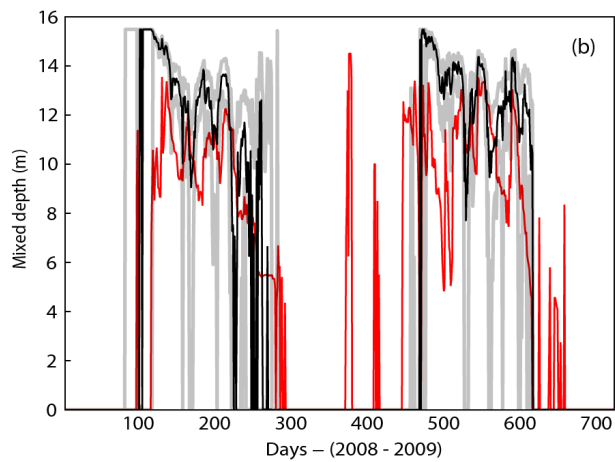
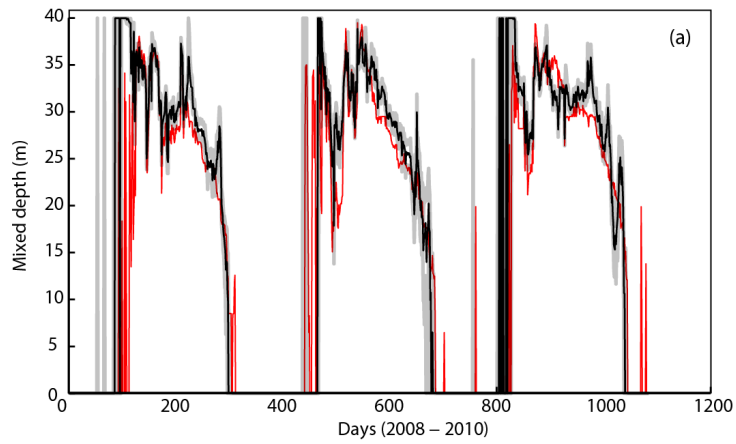
### 392 **3.2.2 Chlorophyll *a* forecasts**

393 For all lake-years, multiple runs of the EM50 Forecasts gave inconsistent simulations and a higher EM size  
394 was required. Forecasts for Windermere tended towards stability between the EM100 and EM200  
395 scenarios (Fig. 3), which is an ensemble size consistent with previous work with relatively complex models  
396 (e.g. Evensen, 1994 and Allen *et al.*, 2003). For Esthwaite Water, however, a higher ensemble size  
397 appeared to be required with a size of around 400 giving consistent simulations (Fig. 4). Subsequently, in  
398 the following, results presented for Windermere and Esthwaite Water are associated with the EM200 and  
399 EM400 scenarios respectively. In all cases, the manually “optimised” variance inflation factor was kept  
400 consistent for all lake years at a value of 1.1.

401

402

403



404

405

406

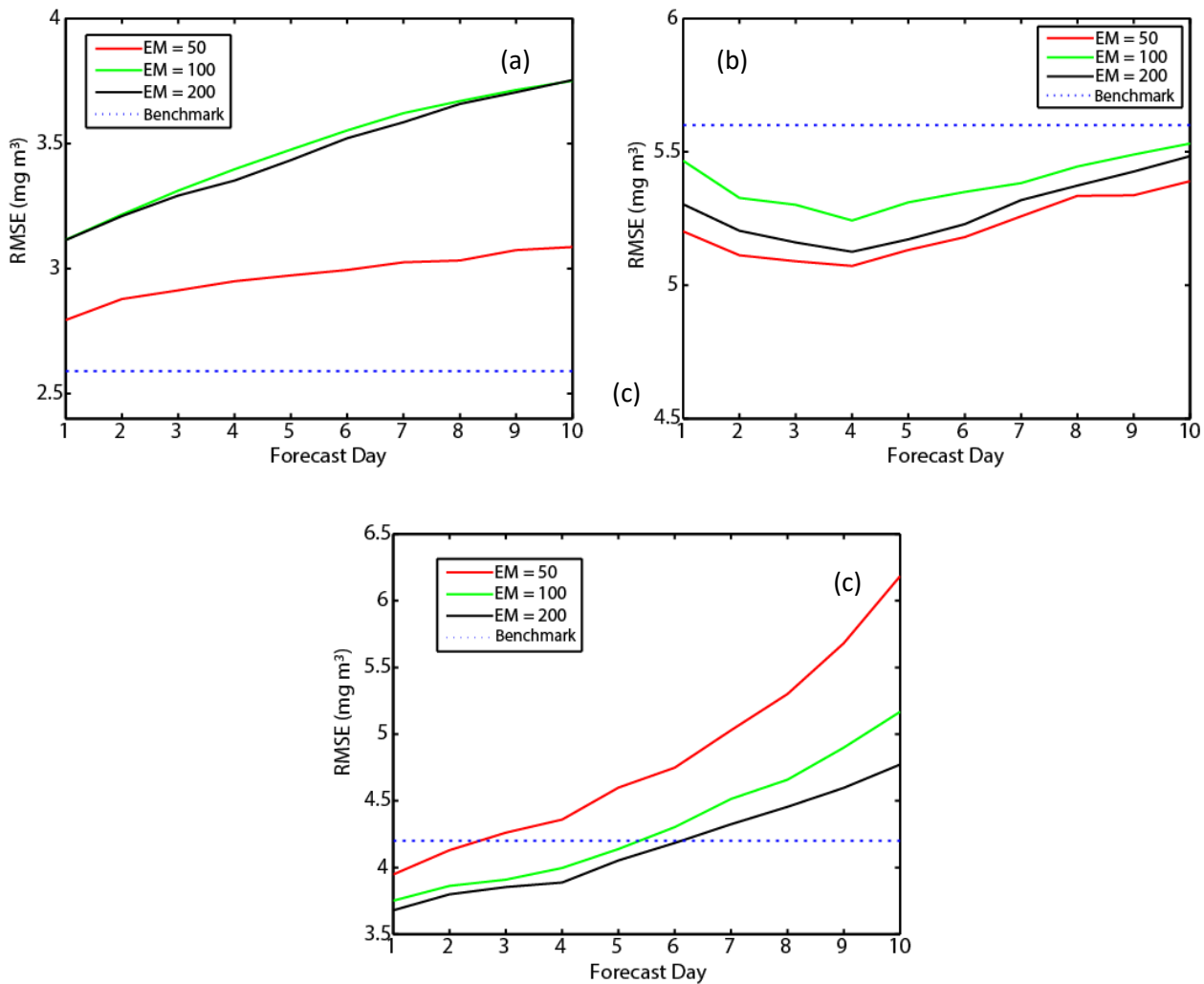
407

408

*Figure 2. Simulated and measured epilimnetic depth. Results shown for (a) Windermere 2008-2010 10-day-ahead, (b) Esthwaite Water 2008 and 2009 10-day-ahead and (c) Esthwaite Water 2008 and 2009 3-day-ahead: “observed” epilimnetic depth (red line), 50<sup>th</sup> percentile of the ensemble of simulated epilimnetic depth (black line) and 5<sup>th</sup> and 95<sup>th</sup> percentiles (grey lines).*

409

410



411

412 *Figure 3. Chlorophyll a forecast skill for the differing ensemble size scenarios. Results are shown*

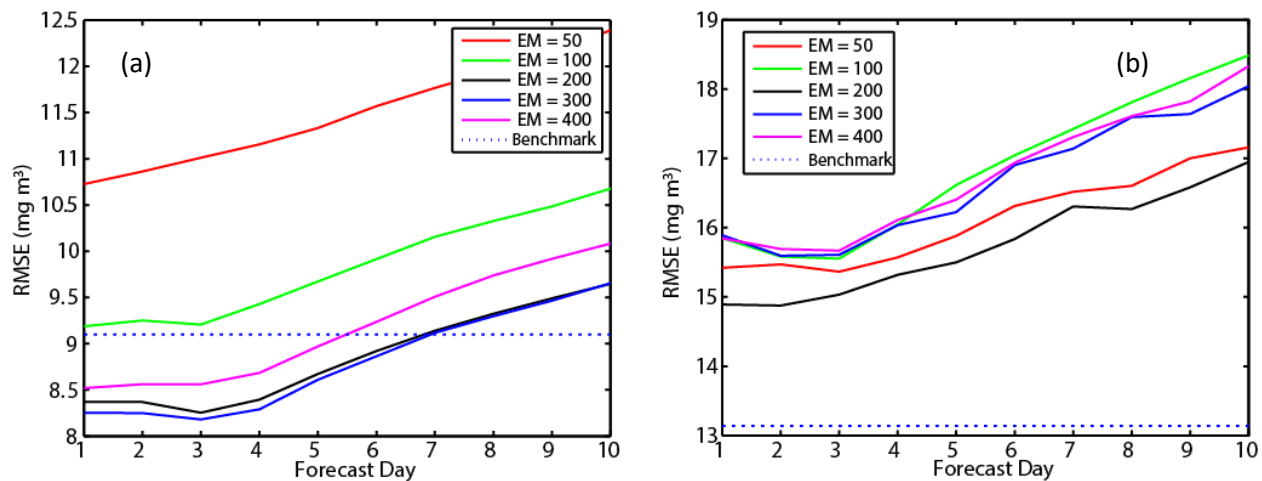
413 *for (a) Windermere 2008, (b) Windermere 2009 and (c) Windermere 2010, compared to the*

414 *benchmark persistence forecast. Note that lower ensemble sizes can give “randomly” better*

415 *forecast performance (e.g. EM = 50 in pane (a))*

416

417



418

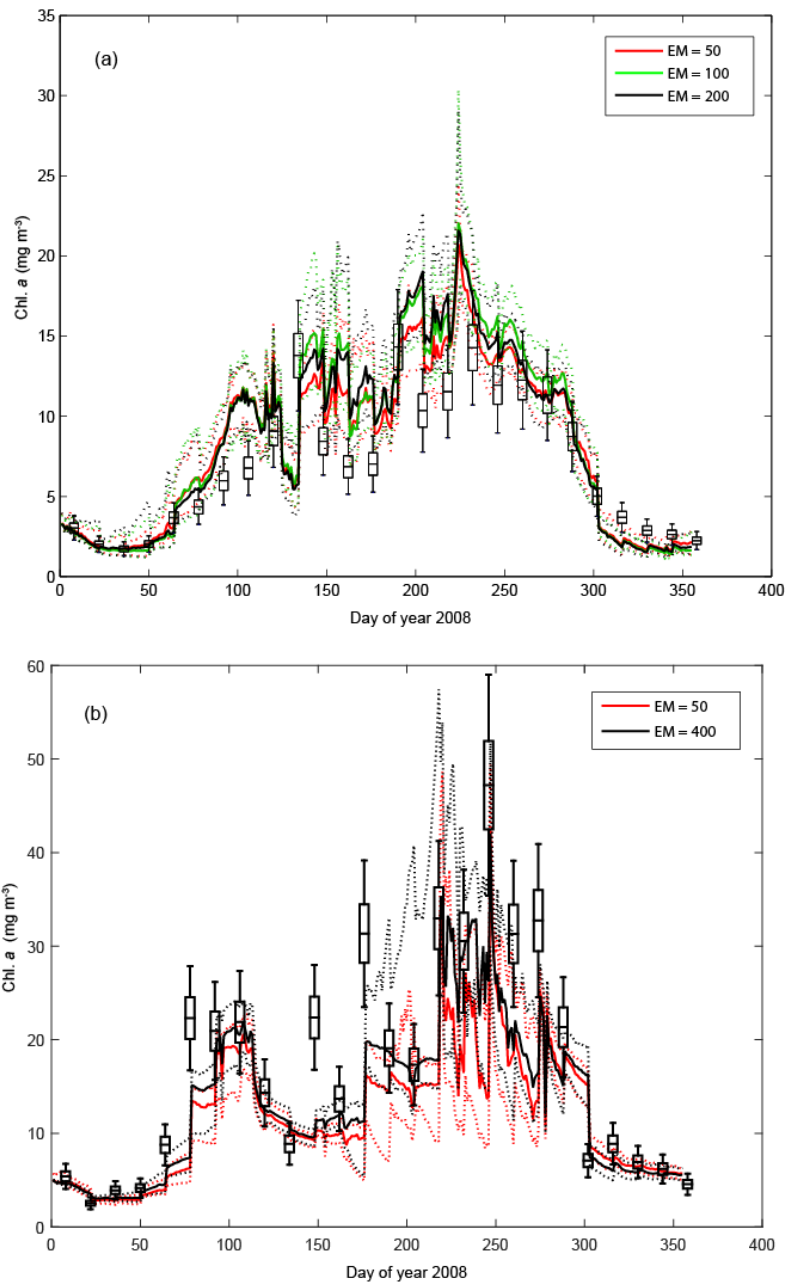
419 *Figure 4. Chlorophyll a forecast skill for the differing ensemble size scenarios. Results are shown*  
420 *for (a) Esthwaite Water 2008 and (b) Esthwaite Water 2009, compared to the benchmark*  
421 *persistence forecast.*

422 Although forecast simulations for Windermere appear to be relatively good visually (e.g. see Fig. 5 below),  
423 they were not always an improvement on the persistence forecasts (Fig. 3). For 2008, the persistence  
424 forecast was better than simulated forecasts for all lead times. Conversely, simulated forecasts were  
425 better than the persistence forecasts for all lead times for 2009. A lead time of approximately 6 days or  
426 less was an improvement on the persistence forecast for 2010 simulations.

427 For Esthwaite Water, forecasts simulations were not as good as those for Windermere (Fig. 5), which is  
428 consistent with previous work using PROTECH for these lakes (Page *et al.*, 2017). The forecasts for 2008  
429 were, however, still better than the persistence forecast out to about 5 days ahead (Fig. 4a), but were  
430 always worse than the persistence forecast for 2009 (Fig. 4b). The poorer fits for Esthwaite Water are  
431 likely to be a result of the complex uncertainties associated with the timing and magnitude of SRP inputs  
432 as well as the poorer simulation of epilimnetic depth reported above. In Esthwaite Water, during the  
433 period where P limitation dominates phytoplankton growth, it is very difficult to represent SRP fluxes

434 appropriately, even when a representation of sediment-derived SRP fluxes was included (the addition of  
435 representation of sediment-derived SRP did not improve forecasts owing to interaction between sources  
436 of P: this work is not reported here). The difficulties associated with representing SRP fluxes was helped  
437 to some degree by the DA, but remain problematic during times where very low concentrations are  
438 present in the epilimnion; at these times, the correlations within the Kalman gain matrix would need to  
439 be very well-represented to provide appropriate updates to both epilimnetic SRP concentrations and SRP  
440 fluxes simultaneously. The difficulties associated with these updates are compounded by the relatively  
441 low frequency of assimilation timesteps. Subsequently, even with relatively large ensemble sizes, the  
442 correlations within the Kalman gain matrix have the potential to be spurious. This is not unexpected as  
443 the lake system is highly dynamic and non-linear and, perhaps most importantly, the relationships  
444 between the states (and parameters in some cases) are not always consistent (e.g. when the nutrient  
445 states are not limiting they may have no relationship with the phytoplankton state). The temporal  
446 evolution of the nutrient parameter values (modified within the DA scheme) that change SRP fluxes were  
447 consistent with these uncertainties and did not show any consistent structure. Given these difficulties,  
448 assimilation of higher resolution nutrient observations may be one of the most important for improving  
449 forecasts. Conversely, for both Windermere and Esthwaite Water, improvement of forecasts was made  
450 by the modification of the background light extinction parameter,  $\epsilon_b$ , within the DA scheme: its evolution  
451 over the simulation periods was relatively consistent for each of the years considered (Fig. 6) and reflects  
452 known simulation artefacts previously reported (Page *et al.*, 2017).





453

454

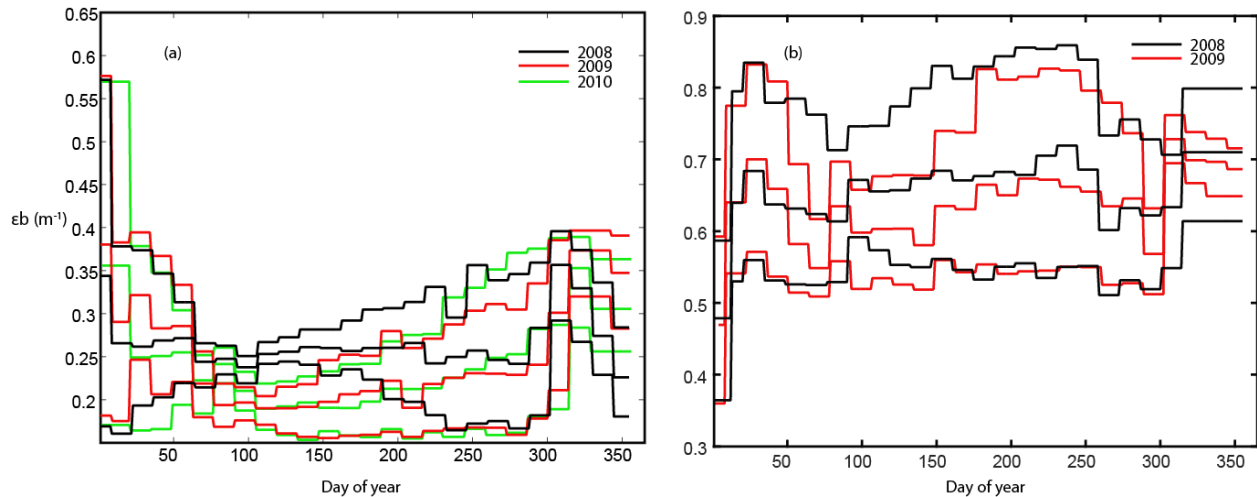
455

456

457

458

*Figure 5. Measured and forecast phytoplankton chlorophyll a in the two lakes during 2008. Results show concatenated forecasts for: (a) 10-day-ahead for Windermere 2008 for ensemble member sizes (EM) of 50, 100 and 200; (b) 5-day-ahead for Esthwaite Water 2008 for ensemble member sizes (EM) of 50 and 400. Solid lines are 50<sup>th</sup> percentile of ensemble and dotted lines are 5<sup>th</sup> and 95<sup>th</sup> percentiles.*



459

460

461

462

463

*Figure 6. The evolution of the background light extinction coefficient parameter ( $\epsilon_b$ ). Results are shown for (a) Windermere 2008, 2009 and 2010 and (b) Esthwaite Water 2008 and 2009. The three lines in each colour are the 5<sup>th</sup>, 50<sup>th</sup> and 95<sup>th</sup> percentiles of the EM200 and EM400 ensembles respectively.*

464

### 3.2.3 Forecasting phytoplankton community structure

465

466

467

468

469

470

Forecasts of species representing the phytoplankton community structure were made without direct constraint within the DA scheme. Simulations were, however, indirectly constrained by the assimilation of mixed depth, chlorophyll *a* and nutrients and hence are reliant on the ability of PROTECH simulations to represent phytoplankton community structure where abiotic conditions for phytoplankton growth are simulated adequately. They are also reliant on whether or not the algal species chosen to represent the community are adequate (Elliott, 2010, 2012; Page *et al.*, 2017).

471

472

473

474

Forecasts of community structure are assessed here using simulations of R- and CS-types functional groups as they dominate our study lakes. Observations to which they are compared are estimated using “counts” of algal species classified into the same functional groups. These “count” data are associated with significant uncertainty in terms of the absolute biovolume of each species (and hence functional

475 type) because of errors, which are difficult to quantify, associated with sample heterogeneity, counter  
476 fatigue and between-counter variation (Thackeray et al., 2012) as well as uncertainty associated with  
477 conversion from sample “counts” to biovolume and subsequently to chlorophyll *a*. Accordingly, we used  
478 the relative abundance of each functional type for each observation timestep to partition the observed  
479 chlorophyll *a* concentration. Given these uncertainties, we estimated the sampling/analytical error to be  
480 +/- 25% and the overall error to be +/- 50% in accordance with Page *et al.* (2017).

481 A comparison of the uncertain observations of R- and CS- functional types are presented in Fig. 7 where  
482 it can be seen that for most lake-years the overall pattern of the simulations are consistent with the  
483 observations. There are some periods where the simulations are not consistent, which are associated  
484 primarily with the period of transition between the early blooms of R-type species and succession by CS-  
485 types (approximately between days 100 and 200). This pattern can clearly be seen for Windemere 2008  
486 and 2009 (Figs. 7a and 7d) and is most likely associated with inadequate representation of nutrient fluxes  
487 and subsequent periods of nutrient limitation (Page *et al.*, 2017). There are also some periods where the  
488 overly rapid mixing simulated by the epilimnetic depth model (as discussed above) made it difficult to  
489 simulate the relatively high observed biomass: this is particularly evident for CS-species in Esthwaite  
490 Water 2008 (Fig. 7k) and R-species in Esthwaite Water 2009 (Fig. 7l); these inconsistencies are a direct  
491 result of the spurious deep mixing events simulated around days 220 and 250 for 2008 and 2009  
492 respectively (see Fig. 2 b and c) and strengthen the requirement to improve the epilimnetic depth model  
493 as discussed above.

#### 494 **3.2.4 Forecasting cyanobacteria**

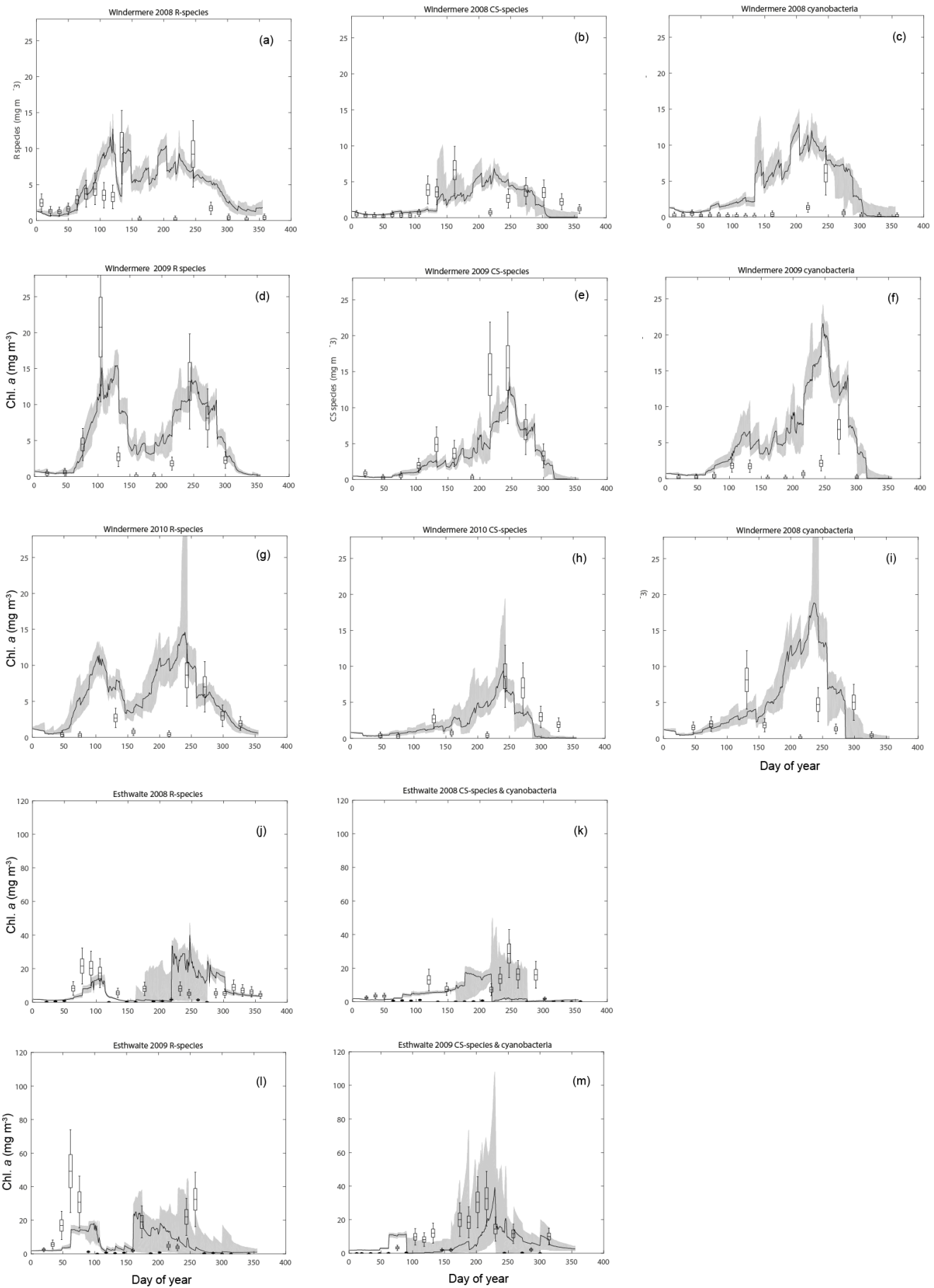
495 Observations of Cyanobacteria are estimated in the same way as functional species types discussed in the  
496 previous section and are associated with similar uncertainty (see Fig. 7). As PROTECH simulates the  
497 functional algal community using the dynamics of a number of selected individual species, the philosophy

498 behind this method means that the forecasts of individual species are not as robust as those for functional  
499 community structure and are hence more uncertain. This is the case for forecasts of cyanobacteria where  
500 they are represented by more than one functional type: e.g. for Windermere cyanobacteria are  
501 represented by *Planktothrix*, an R-type species, together with *Aphanizomenon flos-aquae* and  
502 *dolichospermum* which are CS-type species (see Table Supp. 2). In this situation, the interchangeability of  
503 species with similar functional behaviour, but which have differing species traits, requires additional  
504 interpretation for forecasts of cyanobacteria to be made. For example, the simulations of the R-species  
505 *Planktothrix* for all lake-years for Windermere result in overestimations of cyanobacteria concentrations  
506 for the periods where *Planktothrix* proliferates (approximately between days 150 and 275: Figs. 7c, 7f &  
507 7i). Cyanobacteria forecasts, made for this study, are also a spatial average for each lake, constrained  
508 using data collected at one point; they therefore do not necessarily correspond with the risk from  
509 cyanobacterial blooms where significant spatial heterogeneity exists, as can be the case for wind-blown  
510 cyanobacterial species (e.g. George and Heaney, 1978). Extending point forecasts to spatial forecasts for  
511 species that have these characteristics is hence an additional challenge. However, forecasts may be  
512 presented as probabilistic or possibilistic risk estimates, such as the likelihood of a cyanobacterial  
513 concentration of greater than a given critical threshold: this will be the focus of further research.

#### 514 **4 Conclusions**

515 We rigorously tested the ability of the phytoplankton community model PROTECH to make forecasts of  
516 phytoplankton community structure within a data assimilation scheme using the Ensemble Kalman Filter.  
517 Some forecasting success was shown for chlorophyll *a*, but not all forecasts were better than a persistence  
518 forecast. The results typically indicated a reduction in chlorophyll *a* forecast skill with length of forecasting  
519 period with forecasts for up to four or five days showing greater promise than those for longer time-scales.  
520 Associated forecasts of phytoplankton community composition, represented by functional algal types,

521 were broadly consistent with observations. Translation of forecasts of functional algal types to forecasts  
522 of cyanobacteria are challenging because of functional similarities between species which may or may not  
523 be cyanobacteria. Improvements in forecasts are likely to come from higher frequency observations for  
524 both chlorophyll *a* and nutrient concentrations. - While higher frequency observations for these variables  
525 should help improve forecasts, they will also simultaneously improve the persistence forecast. It,  
526 therefore, remains to be seen whether or not a modelled forecast driven with improved observations  
527 would provide a significant improvement over the associated persistence forecast and the potential to  
528 forecast algal blooms in this type of lake.



530 *Figure 7. Concatenated five-day ahead forecasts of R-species, CS-species and cyanobacteria*  
531 *concentration for all lake years; black line is 50<sup>th</sup> percentile and grey shaded area represents the*  
532 *5<sup>th</sup> and 95<sup>th</sup> percentiles of the ensemble: EM200 and EM400 for Windermere and Esthwaite*  
533 *respectively. The box and whisker symbols represent the analytical uncertainty and the total*  
534 *uncertainty estimated by the project team. Note that 5-day ahead forecasts are presented as*  
535 *approximately this lead time provided the most consistently acceptable results.*

536

### 537 **Acknowledgements**

538 This work was supported by the Natural Environment Research Council projects: the United Kingdom  
539 Lake Ecological Observatory Network (UKLEON; grant number NE/I007407/1) and the Consortium on Risk  
540 in the Environment: Diagnostic, Integration, Benchmarking, Learning and Elicitation (CREDIBLE; grant  
541 number NE/J017299/1); We would like to thank the ECMWF for the historic meteorological forecasts and  
542 Mr Bernard Tebay for collecting the meteorological data at Ambleside.

### 543 **References**

544 Anderson, J.L. (2007). An adaptive covariance inflation error correction algorithm for ensemble filters.  
545 Tellus, 59A, 210–224

546 Anderson, J.L., Anderson, S.L. (1999). A Monte Carlo implementation of the nonlinear filtering problem to  
547 produce ensemble assimilations and forecasts. Mon. Wea. Rev. 127, 2741–2758.

548 J. Allen, M. Eknes, G. Evensen. (2003). An Ensemble Kalman Filter with a complex marine ecosystem  
549 model: hindcasting phytoplankton in the Cretan Sea. Ann. Geophys., 21, 399–411.

550 Bennion, H., Monteith, D. and Appleby, P. (2000). Temporal and geographical variation in lake trophic  
551 status in the English Lake District: evidence from (sub)fossil diatoms and aquatic macrophytes. *Freshwater*  
552 *Biol.*, 45(4), 1365-2427, doi: 10.1046/j.1365-2427.2000.00626.x

553 Buizza, R., Milleer, M. and Palmer, T. N. (1999). Stochastic representation of model uncertainties in the  
554 ECMWF ensemble prediction system. *Q. J. Roy. Meteor. Soc.* 125 (560), 2887–2908.  
555 doi:10.1002/qj.49712556006.

556 Carmichael, W.W. (1992). A status report on planktonic cyanobacteria (blue-green algae) and their toxins.  
557 EPA/600/R-92-079, Environmental Monitoring Systems Laboratory, Office of Research and Development,  
558 U.S. Environmental Protection Agency, Cincinnati, OH. 141 pp.

559 Dong X., Bennion H., Maberly S.C., Sayer C.D., Simpson G.L. and Battarbee R.W. (2012). Nutrients provide  
560 a stronger control than climate on diatom communities in Esthwaite Water Water: Evidence from  
561 monitoring and palaeolimnological records over the past 60 years. *Freshwater Biol.*, 57, 2044-2056.

562 Elliott, J.A. (2012) Predicting the impact of changing nutrient load and temperature on the phytoplankton  
563 of England's largest lake, Windermere. *Freshwater Biol.*, 57, 400-413.

564 Elliott, J.A. (2010). The seasonal sensitivity of Cyanobacteria and other phytoplankton to changes in  
565 flushing rate and water temperature. *Glob. Change Biol.*, 16, 864-876.

566 Evensen, G. (1994). Sequential data assimilation with a non-linear quasigeostrophic model using Monte  
567 Carlo methods to forecast error statistics. *J. Geophys. Res.*, 99, 10 143–10 162.

568 Evensen, G. (2009). The ensemble Kalman filter for combined state and parameter estimation. *IEEE*  
569 *Control Sys.*, 29 (3), pp. 83-104, 2009.



570 George, D. G. and Heaney, S. I. (1978). Factors influencing the spatial distribution of phytoplankton in a  
571 small productive lake. *J. Ecol.*, 66(1), 133-155.

572 Hall, G.H., Maberly, S.C., Reynolds, C.S., Winfield, I.J., James, B.J., Parker, J.E., Dent, M.M., Fletcher, J.M.,  
573 Simon, B.M. and Smith, E. (2000). Feasibility study on the restoration of three Cumbrian lakes. Centre for  
574 Ecology and Hydrology Windermere, Ambleside, UK. 82 pp.

575 Heany, S.I., Corry, J. E. and Lishman, J. P. (1992). Changes of water quality and sediment phosphorus of a  
576 small productive lake following decreased phosphorus loading. Centre for Ecology and Hydrology  
577 Windermere, Ambleside, UK. 14 pp.

578 Ho, J. C. and Michalak, A. M. (2015). Challenges in tracking harmful algal blooms: A synthesis of evidence  
579 from Lake Erie. *J. Great Lakes Res.*, 41(2), 317-325. doi.org/10.1016/j.jglr.2015.01.001

580 Huang, J., Gao, J., Liu, J. and Zhang, Y. (2013). State and parameter update of a hydrodynamic-  
581 phytoplankton model using ensemble Kalman filter, *Ecol. Model.*, 263 (10), 81-91.  
582 <https://doi.org/10.1016/j.ecolmodel.2013.04.022>

583 Kim, K., Park, M., Min, J., Ryu, I., Kang, M., and Park, L. (2014). Simulation of algal bloom dynamics in a  
584 river with the ensemble Kalman filter. *J. Hydrol.*, 519(D), 2810–2821.  
585 <https://doi.org/10.1016/j.jhydrol.2014.09.073>

586 Luo, Y., Ogle, K., Tucker, C., Fei, S., Gao, C., LaDeau, S., Clark, J. S. and Schimel, D. S. (2011). Ecological  
587 forecasting and data assimilation in a data-rich era. *Ecol. Appl.*, 21, 1429–1442. doi:10.1890/09-1275.1

588 Maberly, S.C., De Ville, M.M., Thackeray, S.J., Feuchtmayr, H., Fletcher, J.M., James, J.B., Kelly, J.L., Vincent,  
589 C.D., Winfield, I.J., Newton, A., Atkinson, D., Croft, A., Drew, H., Saag, M., Taylor, S., Titterton, H. (2011).  
590 A survey of the lakes of the English Lake District: The Lakes Tour 2010. NERC/Centre for Ecology and  
591 Hydrology, 137pp. (CEH Project Number. Report to: Environment Agency, North West Region and Lake

592 District National Park Authority: downloaded Jan 2015 from  
593 <http://nora.nerc.ac.uk/14563/2/N014563CR.pdf>

594 Mackay E. M., Folkard A. M. and Jones I.D. (2014). Interannual variations in atmospheric forcing determine  
595 trajectories of hypolimnetic soluble reactive phosphorus supply in a eutrophic lake. *Freshwater Biol.*, 59,  
596 1646–1658.

597 Madgwick G., Jones I.D., Thackeray S.J., Elliott J.A. and Miller H.J. (2006). Phytoplankton communities and  
598 antecedent conditions: high resolution sampling in Esthwaite Water. *Freshwater Biol.*, 51, 1798–  
599 1810.

600 Marcé, R., George, G., Buscarinu, P., Deidda, M, Dunalska, J., de Eyto, E., Flaim, G., Grossart, H.,  
601 Istvanovics, V., Lenhardt, M., Moreno-Ostos, E., Obrador, B., Ostrovsky, I., Pierson, D. C., Potužák, J.,  
602 Poikane, S., Rinke, K., Rodríguez-Mozaz, S., Staehr, P. A., Šumberová, K., Waajen, G., Weyhenmeyer, G. A.,  
603 Weathers, K. C., Zion, M., Ibelings, B.W. and Jennings, E. (2016). Automatic High Frequency Monitoring  
604 for Improved Lake and Reservoir Management. *Environ. Sci. Technol.* 50 (20), 10780-10794. DOI:  
605 10.1021/acs.est.6b01604

606 Michalak, A., M. (2016). Study role of climate change in extreme threats to water quality. *Nature* 535, 349-  
607 350.

608 Metcalf, J.S. and Codd, G.A. (2009). Cyanobacteria, neurotoxins and water resources: are there  
609 implications for human neurodegenerative disease? *Amyotrophic Lateral Sclerosis* 10, suppl. 2, 74-78  
610 (2009).

611 Moradkhani, H., S. Sorooshian, H. V. Gupta, and P. R. Hauser (2005). Dual state-parameter estimation of  
612 hydrological models using ensemble Kalman filter, *Adv. Water Resour.*, 28, 135 – 147.

613 Ollinaho, P., Lock, S.-J., Leutbecher, M., Bechtold, P., Beljaars, A., Bozzo, A., Forbes, R. M., Haiden, T.,  
614 Hogan, R. J. and Sandu, I. (2017), Towards process-level representation of model uncertainties:  
615 stochastically perturbed parametrizations in the ECMWF ensemble. *Q.J.R. Meteorol. Soc.*, 143: 408–422.  
616 doi:10.1002/qj.2931

617 Page et al., (2017). Constraining uncertainty and process-representation in an algal community lake model  
618 using high frequency in-lake observations. *Ecol. Model.*:  
619 <http://www.sciencedirect.com/science/article/pii/S0304380017301345>

620 Paerl, H.W. and Huisman, J. (2008). Blooms like it hot. *Science*, 4, 320(5872), 57-8. doi:  
621 10.1126/science.1155398. DOI: 10.1126/science.1155398

622 Pretty, J. N., Mason, C. F., Nedwell, D. B., Hine, R. E., Leaf, S., and Dils, R. (2003). Environmental Costs of  
623 Freshwater Eutrophication in England and Wales. *Environ. Sci. Technol.*, 37(2), 201-208.

624 Read J.S., Hamilton, D.P., Jones, I.D., Muraoka, K., Winslow, L.A. , Kroiss, R. , Wu, C.H. & Gaiser. E. (2011).  
625 Derivation of lake mixing and stratification indices from high-resolution lake buoy data. *Environ. Modell.*  
626 *Softw.* 26, 1325-1336.

627 Ramsbottom A.E. (1976). Depth Charts of the Cumbrian Lakes. Freshwater Biological Association Scientific  
628 Publication No. 33, Ambleside, UK.

629 Recknagel, F., Ostrovsky, I. and Cao, H. (2014). Model ensemble for the simulation of plankton community  
630 dynamics of Lake Kinneret (Israel) induced from in situ predictor variables by evolutionary computation.  
631 *Environ. Modell. Softw.*, 61, 380-392. <https://doi.org/10.1016/j.envsoft.2014.03.014>.

632 Reynolds C.S. (1988). Functional morphology and the adaptive strategies of freshwater phytoplankton. In:  
633 Growth and Reproductive strategies of Freshwater Phytoplankton (Ed. C.D. Sandgren), pp. 388–433.  
634 Cambridge, University Press, New York.

635 Reynolds C.S., Irish A.E. and Elliott J.A. (2001). The ecological basis for simulating phytoplankton responses  
636 to environmental change (PROTECH). *Ecol. Model.*, 140, 271–291.

637 Rigosi, A., Carey, C.C., Ibelings, B. W. and Brookes, J. D. (2014). The interaction between climate warming  
638 and eutrophication to promote cyanobacteria is dependent on trophic state and varies among taxa.  
639 *Limnol. Oceanogr.* 59(1), 2014, 99–114. doi:10.4319/lo.2014.59.01.0099

640 Rowe, M. D., Anderson, E. J. , Wynne, T. T., Stumpf, R. P., Fanslow, D. L., Kijanka, K., Vanderploeg, H. A.  
641 Strickler, J. R. and Davis, T. W. (2016). Vertical distribution of buoyant *Microcystis* blooms in a Lagrangian  
642 particle tracking model for short-term forecasts in Lake Erie. *J. Geophys. Res.: Oceans.* 121, 5296-5314.  
643 doi:10.1002/2016JC011720.

644 Smith, V.H., (2003). Eutrophication of Freshwater and Coastal Marine Ecosystems: A Global Problem.  
645 *Environ. Sci. & Pollut. Res.* 10 (2) 126-39.

646 Stumpf, R. P., Tomlinson, M. C., Calkins, J. A., Kirkpatrick, B., Fisher, K., Nierenberg, K., Currier, R. and  
647 Wynne, T. T. (2009). Skill assessment for an operational algal bloom forecast system. *Journal of Marine*  
648 *Systems.* 76(1): 151-161.

649 Taylor, C.J., Pedregal, D.J., Young, P.C. and Tych, W., (2007). Environmental time series analysis and  
650 forecasting with the Captain toolbox, *Environ. Modell. Softw.*, 22: 797-814.

651 World Health Organization (1999). Toxic cyanobacteria in water: a guide to their public health  
652 consequences, monitoring and management. I. Chorus and J. Bartram (Eds.). E & FN Spon, London, UK  
653 (1999).

654

655 Xiao X, Sogge H, Lagesen K, Tooming-Klunderud A, Jakobsen KS, Rohrlack T (2014). Use of High Throughput  
656 Sequencing and Light Microscopy Show Contrasting Results in a Study of Phytoplankton Occurrence in a  
657 Freshwater Environment. PLoS ONE, 9(8), 1-9. doi:10.1371/journal.pone.0106510

658 Xiao, X., He, J., Huang, H., Miller, T. R., Christakos, G., Reichwaldt, E., S., Ghadouani, A., Lin, S., Xu, X. and  
659 Shi, J. (2017). A novel single-parameter approach for forecasting algal blooms. Water Res., 108, 222-231.  
660 <https://doi.org/10.1016/j.watres.2016.10.076>

661 Young, P.C. and Beven, K.J. 1994. Data-based mechanistic modelling and the rainfall-flow non-linearity.  
662 Environmetrics. 5, 3, p. 335-363.

663 Ye, L., Cai, Q., Zhang, M. and Tan, L. (2014). Real-time observation, early warning and forecasting  
664 phytoplankton blooms by integrating in situ automated online sondes and hybrid evolutionary algorithms.  
665 Ecological Informatics, 22, 44–51.

666 Young, P.C., 2015. Refined Instrumental Variable Estimation: Maximum Likelihood Optimization of a  
667 Unified Box-Jenkins Model. Automatica, 52, 35–46.

## 668 **Supplementary information**

### 669 **Supp. 1 Transfer Function models for forecasted inputs**

670 The epilimnetic depth model requires forecasts of epilimnetic temperature, river in/outflows and river  
671 temperature. Each TF model that provides these forecasts was identified (as outlined above) using the  
672 available timeseries data. The epilimnetic temperature ( $T_e$ ) at day  $t$  is given by:

673

674

$$675 \quad T_{e(t)} = -a \cdot T_{e(t-1)} + b_1 \cdot T_a(t) + b_2 \cdot R_{sw(t)} + b_3 \cdot \frac{1}{D_{e(t-1)}} + b_4 \cdot (W_{s(t-1)})^3$$

676

677

678 Where,  $T_a$  is the air temperature,  $R_{sw}$  is SW radiation,  $D_e$  is epilimnetic depth and  $W_s$  is the wind speed.

679 The model coefficients are denoted  $a$ ,  $b1$ ,  $b2$  and  $b3$  (see Table Supp. 1 for values). One model for each

680 lake was identified from the available data (2008 to 2010 for Windermere and 2004 to 2009 for Esthwaite

681 Water).

682 The lake in/outflow TF model was identified as a 1<sup>st</sup> order model with a nonlinear rainfall filter (see Young

683 and Beven, 1994) and took the form:

684

$$685 \quad Q_{r(t)} = -a \cdot Q_{r(t-1)} + b \cdot P_{(t)} \cdot Q_{r(t-1)}^\beta$$

686

687

688 where  $Q_r$  is the river in/outflow,  $P$  is precipitation and  $a$ ,  $b1$  are TF model coefficients where  $\beta$  is the

689 nonlinear rainfall filter parameter. The model for Windermere was identified using Rainfall data from

690 Ambleside and flow data from the Environment agency Gauge at Newby Bridge for the years 2008 to 2010

691 (National River Flow Archive: <http://www.ceh.ac.uk/data/nrfa/>).

692 River temperature ( $T_Q$ ) was estimated using observed data from Troutbeck (Windermere) for the years

693 1997 to 2006:

694

$$695 \quad T_{Q(t)} = -a \cdot T_{Q(t-1)} + b \cdot T_{a(t)}$$

696

697

698

699 **Table Supp. 1 Transfer Function parameters and goodness of fit (W = Windermere; E = Esthwaite Water)**

	a		b1 ( $\beta$ )		b2		b3		b4		$\tau$		$R_r^2$	
	W	E	W	E	W	E	W	E	W	E	W	E	W	E
Lake Surface Temperature ( $T_s$ )	-0.9449	-0.899	0.055	0.093	0.0008	0.0025	0.0011	0.0022	-0.0007	-0.0012	[0,0,0,0]	[0,1,1,0]	0.97	0.98
River in/outflow ( $Q_r$ )	-0.7717	-0.829	11.141 (0.2)	0.022 (0.3)			-	-			1	0	0.92	0.86
River Temperature ( $T_d$ )	-0.900	-0.900	0.1005	0.1005	-	-	-	-	-	-	0	0	0.87	0.87

700

701 **Table Supp. 2. Species used to represent algal communities. Functional algal types and an indication of**  
 702 **classification as cyanobacteria given are in parenthesis: functional types follow Reynolds (1988).**

Windermere	Esthwaite Water Water
<i>Aphanizomenon flos-aquae</i> (CS; Cyano)	<i>Asterionella</i> (R)
<i>Aulacoseira</i> (R)	<i>Aulacoseira</i> - 2008 (R); <i>Fragilaria crotonensis</i> -(2009 (R)
<i>Asterionella</i> (R)	<i>Aphanizomenon flos-aquae</i> (CS; Cyano)
<i>Cryptomonas</i> (CSR)	<i>Aphanothece clathrata</i> (CS; Cyano)
<i>Dolichospermum</i> (CS; Cyano)	<i>Cryptomonas</i> (CSR)
<i>Monoraphidium</i> (CS)	<i>Dictyosphaerium pulchellum</i> (R)
<i>Paulschulzia tenera</i> (S)	<i>Dolichospermum</i> (CS; Cyano)
<i>Planktothrix</i> (R; Cyano)	<i>Eudorina</i> (S)

703