

# 1 Using feature-based verification methods to explore the spatial and 2 temporal characteristics of the 2019 Chlorophyll-*a* bloom season in a 3 model of the European North-West Shelf

4 Marion Mittermaier<sup>1</sup>, Rachel North<sup>1</sup>, Jan Maksymczuk<sup>2</sup>, Christine Pequignet<sup>2</sup>, David Ford<sup>2</sup>

5 <sup>1</sup>Verification, Impacts and Post-Processing, Weather Science, Met Office, Exeter, EX1 3PB, United Kingdom

6 <sup>2</sup>Ocean Forecasting Research & Development, Weather Science, Met Office, Exeter, EX1 3PB, United Kingdom

7  
8 *Correspondence to:* Marion Mittermaier (marion.mittermaier@metoffice.gov.uk)

## 9 **Abstract.**

10 Two feature-based verification methods, used for atmospheric model applications, have been applied to  
11 compare pre-operational analyses of Chlorophyll-*a* (Chl-*a*) concentrations from the Met Office Atlantic  
12 Margin Model at 7 km resolution (AMM7v11) for the North West European Shelf Seas with a gridded  
13 satellite-derived Chl-*a* concentrations product from the Copernicus Marine Environment Monitoring  
14 Service (CMEMS) catalogue. Chl-*a* bloom objects were identified using a range of thresholds for the  
15 2019 bloom season (March 1 to 31 July). These bloom objects were analysed as purely spatial features  
16 and as space-time objects, enabling the ability to define the onset, duration and demise of distinct bloom  
17 episodes. Overall, the AMM7v11 analyses were found to be similar to the satellite product. The  
18 AMM7v11 analyses were not always able to represent coastal objects given the coastline definition in a  
19 ~7 km model and sub-grid scale processes. By contrast the AMM7v11 analyses produces more bloom  
20 objects in deeper Atlantic waters, which are not detected by the satellite product. Concentrations in the  
21 AMM7v11 analyses are somewhat higher overall. This bias manifests itself in the size of the  
22 AMM7v11 bloom objects, which tend to be larger than the bloom objects identified in the satellite  
23 product. Based on this analysis the onset of the bloom season is delayed by 26 days in the AMM7v11,  
24 but the season also persists for another month beyond the diagnosed end. Overall, the season was  
25 diagnosed to be 119 days long, based on the AMM7v11 space-time objects, and 117 days from the  
26 satellite product. Geographically the AMM7v11 and satellite product objects do overlap at times, but  
27 further analysis shows they do not necessarily exist in that location at the same time.

## 28 **1 Introduction**

29 The advancements in atmospheric numerical weather prediction (NWP) such as the improvements in  
30 model resolution began to expose the relative weaknesses in so-called traditional verification scores  
31 (such as the root-mean-squared-error for example), which rely on the precise matching in space and  
32 time of the forecast to a suitable observation. These metrics and measures no longer provided adequate  
33 information to quantify forecast performance (e.g. Mass et al. 2002). One key characteristic of high-  
34 resolution forecasts is the apparent detail they provide, but this detail may not be in the right place at the  
35 right time, a phenomenon referred to as the “double penalty effect” (Rossa et al., 2008). Essentially it  
36 means that at any given time the error is counted twice because the forecast occurred where it was not  
37 observed, and it did not occur where it was observed. This realisation created the need within the  
38 atmospheric community for creating more informative yet robust verification methods. As a result, a  
39 multitude of so-called “spatial” verification methods were developed, which essentially provide a  
40 number of ways for accounting for the characteristics of high-resolution forecasts.

41

42 In 2007 a spatial verification method inter-comparison (Gilleland et al., 2009, 2010) was established  
43 with the aim of providing a better collective understanding of what each of the new methods was  
44 designed for, and categorising what type of forecast errors each could quantify. A decade later  
45 Dorninger et al. (2018) revisited this inter-comparison, adding a fifth category so that all spatial  
46 methods fall into one of the following groupings: neighbourhood, scale separation, feature-based,  
47 distance metrics or field deformation.

48

49 The use of spatial verification methods has therefore become commonplace for atmospheric NWP (see  
50 Dorninger et al. (2018) and references within). Neighbourhood-based methods in particular have  
51 become popular due to the relative ease of computation and intuitive interpretation. Recently one such  
52 neighbourhood spatial method was demonstrated as an effective approach for exploring the benefit of  
53 higher resolution ocean forecasts (Crocker et al., 2020). Another class of methods focus on how well  
54 particular features of interest are being forecast. Forecasting specific features of interest is one of the  
55 main reasons for increasing horizontal resolution. Feature-based verification methods, such as the

56 Method for Object-based Diagnostic Evaluation (MODE, Davis et al., 2006) and the time domain  
57 version MODE-TD (Clark et al., 2014) enable an assessment of such features, focusing on the physical  
58 attributes of the features (identified using a threshold) and how they behave at a given point in time, and  
59 evolve over time. These methods require a gridded truth to compare to. Whilst the initial inter-  
60 comparison project was based on analysing precipitation forecasts, over recent years their use has  
61 extended to other variables, provided gridded data sets exist that can be used to compare against (e.g.  
62 Crocker & Mittermaier (2013) considered cloud masks and Mittermaier et al., (2016) considered more  
63 continuous fields in a global NWP model such as upper-level jet cores, surface lows and high pressure  
64 cells using model analyses. Mittermaier & Bullock (2013) detailed the first study to use MODE-TD  
65 prototype tools to analyse the evolution of cloud breaks over the UK using satellite-derived cloud  
66 analyses.

67

68 In the ocean, several processes have strong visual signatures that can be detected by satellite sensors.  
69 For example, mesoscale eddies can be detected from sea surface temperature or sea level anomaly (e.g.  
70 (Chelton et al., 2011, Morrow and Le Traon, 2012, Hausmann and Czaja, 2012). Phytoplankton blooms  
71 are seasonal events which see rapid phytoplankton growth as a result of changing ocean mixing,  
72 temperature and light conditions (Sverdrup, 1953, Winder and Cloern, 2010, Chiswell, 2011)). Blooms  
73 represent an important contribution to the oceanic primary production that is a key process for the  
74 oceanic carbon cycle (Falkowski et al., 1998). Their spatial extent and intensity in the upper ocean make  
75 them visible from space with ocean colour sensors (Gordon et al., 1983, Behrenfeld et al., 2005).  
76 Biogeochemical models coupled to physical models of the ocean provide simulations for the various  
77 parameters that characterise the evolution of a spring bloom. In particular, Chlorophyll-*a* (Chl-*a*)  
78 concentrations provide an index of phytoplankton biomass. Chl-*a* concentration can also be estimated  
79 from spaceborne ocean colour sensors (Antoine et al., 1996).

80

81 Validation of marine biogeochemical models has traditionally relied on simple statistical comparisons  
82 with observation products, often limited to visual inspections (Stow et al., 2009; Hipsey et al., 2020). In  
83 response to this, various papers have outlined and advocated using a hierarchy of statistical techniques

84 (Allen et al., 2007a, 2007b; Stow et al., 2009; Hipsey et al., 2020), multivariate approaches (Allen and  
85 Somerfield, 2009), and novel diagrams (Jolliff et al., 2009). Many of these rely on matching to  
86 observations in space and time, but some studies have started applying feature-based verification  
87 methods. Emergent properties have been assessed in terms of geographical provinces (Vichi et al.,  
88 2011), phenological indices (Anugerahanti et al., 2018), and ecosystem functions (De Mora et al.,  
89 2016). In a previous application of spatial verification methods developed for NWP, Saux Picart et al.,  
90 2012) used a wavelet-based method to compare Chl-*a* concentrations from a model of the European  
91 North West Shelf to an ocean colour product.

92

93 For this paper, both MODE and MODE-TD (or MTD for short) were applied to the latest pre-  
94 operational analysis (at the time) of the Met Office Atlantic Margin Model (AMM7) at 7 km resolution  
95 (O’Dea et al., 2012; Edwards et al., 2012; O’Dea et al., 2017; King et al., 2018) for the European North  
96 West Shelf (NWS), in order to evaluate the spatio-temporal evolution of the bloom season in both  
97 model and observation fields. A traditional verification of the system (e.g. using root-mean-squared-  
98 error and similar metrics) is out of scope of this study and will be presented in a separate publication.  
99 Traditional verification of a previous version, prior to the introduction of ocean colour data assimilation,  
100 was presented by Edwards et al. (2012), who used various metrics and Taylor diagrams (Taylor, 2001)  
101 to compare model analyses to satellite and in-situ observations. Ford et al. (2017) presented further  
102 validation, to understand the skill of the model at representing phytoplankton community structure in  
103 the North Sea. A similar version of the system used in this study, including ocean colour data  
104 assimilation, was assessed in Skákala et al. (2018), who validated both analysis and forecast skill using  
105 traditional methods. The assimilation improved analysis and forecast skill compared with the free-  
106 running model, but when assessed against satellite ocean colour the forecasts were not found to beat  
107 persistence. On the NWS the spring bloom usually begins between February and April, varying across  
108 the domain and interannually (Siegel et al., 2002; Smyth et al., 2014), and lasts until summer. Without  
109 data assimilation the spring bloom in the model typically occurs later than in observations (Skákala et  
110 al., 2018, 2020), a bias which is largely corrected by assimilating ocean colour data.

111

112 In Section 2 the data sets used in the verification process are introduced. Section 3 describes MODE and  
113 MTD. Section 4 contains a selection of results, and their interpretation. Conclusions and  
114 recommendations follow in Section 5.

## 115 **2 Data sets for the 2019 Chl-*a* bloom**

116 As stated in Section 1, feature-based methods such as MODE and MTD require the fields to be  
117 compared to be on the same grid.

### 118 **2.1 Satellite-derived gridded ocean colour products**

119 A cloud-free gridded (space-time interpolated, L4) daily product delivered through the Copernicus  
120 Marine Environment Monitoring Service (CMEMS, Le Traon et al., 2019) catalogue provides Chl-*a*  
121 concentration at ~1 km resolution over the Atlantic (46°W–13°E, 20°N–66°N). The L4 Chl-*a* product is  
122 derived from merging of data from multiple satellite-borne sensors: MODIS-Aqua, VIIRS-N and OLCI-  
123 S3A. The reprocessed (REP) products available nearly 6 months after the measurements  
124 (OCEANCOLOUR\_ATL\_CHL\_L4\_REP\_OBSERVATIONS\_009\_091) are used here as it is the best-  
125 quality gridded product available for comparison. The satellite derived chlorophyll concentration  
126 estimate is an integrated value over optical depth.

127

128 Errors in satellite-derived Chl-*a* can be more than 100% of the observed value (e.g. Moore et al., 2009).  
129 The errors in the L4 Chl-*a* values are often at their largest near the coast, especially near river outflows.  
130 However, in the rest of the domain, smaller values of Chl-*a* mean that even large percentage  
131 observation errors result in errors typically smaller than the difference between model and observations.  
132 As will be shown, the models at 7 km resolution cannot resolve the coasts in the same way as is seen in  
133 the satellite product as some of the coastal Chl-*a* dynamics are sub-grid scale for a 7 km resolution  
134 model.

135

136 For this study the ~1 km resolution L4 satellite product was interpolated onto the AMM7 grid using  
137 standard two-dimensional horizontal cubic interpolation. This coarsening process retained some of the  
138 larger concentrations present in the L4 product.

## 139 **2.2 Model description**

140 Operational modelling of the NWS is performed using the Forecast Ocean Assimilation Model (FOAM)  
141 system. This consists of the NEMO (Nucleus for European Modelling of the Ocean) hydrodynamic  
142 model (Madec et al., 2016; O'Dea et al., 2017), the NEMOVAR data assimilation scheme (Waters et al.,  
143 2015; King et al., 2018), and for the NWS region the European Regional Seas Ecosystem Model  
144 (ERSEM), which provides forecasts for the lower trophic levels of the marine food web (Butenschön et  
145 al., 2016). The version of FOAM used in this study is AMM7v11, using the ~7 km horizontal  
146 resolution domain stretching from 40 °N, 20 °W to 65 °N, 13 °E. Operational forecasts of ocean physics  
147 and biogeochemistry for the NWS are delivered through CMEMS, for a summary of the principles  
148 underlying the service see e.g. Le Traon et al. (2019).

149

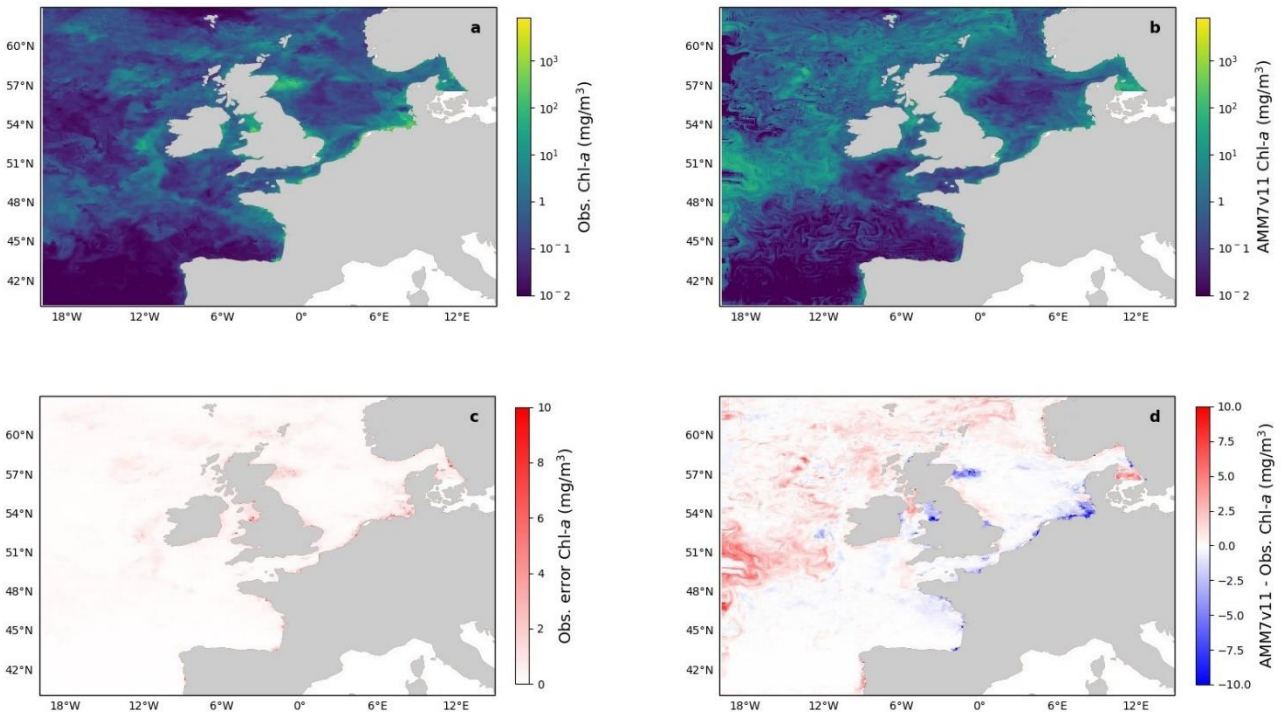
150 AMM7v11 uses the CO6 configuration of NEMO, which is configured for the shallow water of the  
151 shelf sea and is a development of the CO5 configuration described by O'Dea et al. (2017). The ERSEM  
152 version used is v19.04, coupled to NEMO using the Framework for Aquatic Biogeochemical Models  
153 (FABM, Bruggeman and Bolding, 2014). The NEMOVAR version is v6.0, with a 3D-Var method used  
154 to assimilate satellite and in situ sea surface temperature (SST) observations, in situ temperature and  
155 salinity profiles, and altimetry data into NEMO (King et al., 2018), and chlorophyll derived from  
156 satellite ocean colour into ERSEM (Skákala et al., 2018). The introduction of ocean colour assimilation  
157 in AMM7v11 is a major development for the biogeochemistry over previous versions of the system  
158 (Edwards et al., 2012). The satellite ocean colour observations assimilated are from a daily L3 multi-  
159 sensor composite product based on MODIS and VIIRS with resolutions of 1 km for the Atlantic (for  
160 further information see OCEANCOLOUR\_ATL\_CHL\_L3\_NRT\_OBSERVATIONS\_009\_036 on the  
161 CMEMS catalogue).

162

163 In this study daily mean Chl-*a* concentrations for the period of 1 March-31 July 2019 from AMM7v11  
164 were used to illustrate the verification methodology. AMM7v11 entered operational use in December  
165 2020, and the data used here came from a pre-operational run of the system. Note only the analysis of  
166 AMM7v11 (i.e. no corresponding forecasts) was available at the time of the assessment, and the results  
167 presented in this paper show how close the data assimilation draws the model to the observed state.

### 168 **2.3 Visual inspection of data sets**

169 Ideally, Chl-*a* concentration from the model should be integrated over optical depth to be equivalent to  
170 the satellite derived value defined in Section 2.1 (Dutkiewicz et al., 2018). However, this is currently a  
171 non-trivial exercise, and cannot be accurately calculated from offline outputs. Therefore, the commonly  
172 accepted practice is to use the model surface Chl-*a* (Lorenzen, 1970, (Shutler et al., 2011). Here it is  
173 assumed that the difference between surface and optical depth-integrated Chl-*a* is likely to be small in  
174 comparison with the actual model errors.



175

176 **Figure 1 (a) Daily mean L4 multi-sensor observations regrided on the 7 km resolution model grid and (b) AMM7v11**  
 177 **Chl-*a* for 1 June 2019. (c) Error estimates on the multi-sensor L4 Chl-*a* and (d) difference between AMM7v11 and**  
 178 **the L4 product.**

179

180 Figure 1 shows the L4 ocean colour product (a) and AMM7v11 analysis (b) for 1 June 2019 on the top  
 181 row, using the same plotting ranges. The second row shows the difference field that is provided with the  
 182 L4 ocean colour product (c), and the AMM7v11 minus L4 difference field (d). The mean error (bias) is  
 183 generally positive with the AMM7v11 analysis containing higher Chl-*a* concentrations, especially in the  
 184 deeper North Atlantic waters. The exceptions are along the coast where the AMM7v11 analysis is  
 185 deficient, but it should be noted that these are also the zones where some of the largest satellite retrieval  
 186 errors occur and where a 7-km resolution model, with a coarse representation of the coast, does not fully  
 187 represent complex coastal and estuarine processes.

### 188 **3 Method for Object-based Diagnostic Evaluation (MODE) and MODE Time-Domain (MTD)**

#### 189 **3.1. Description of the methods**



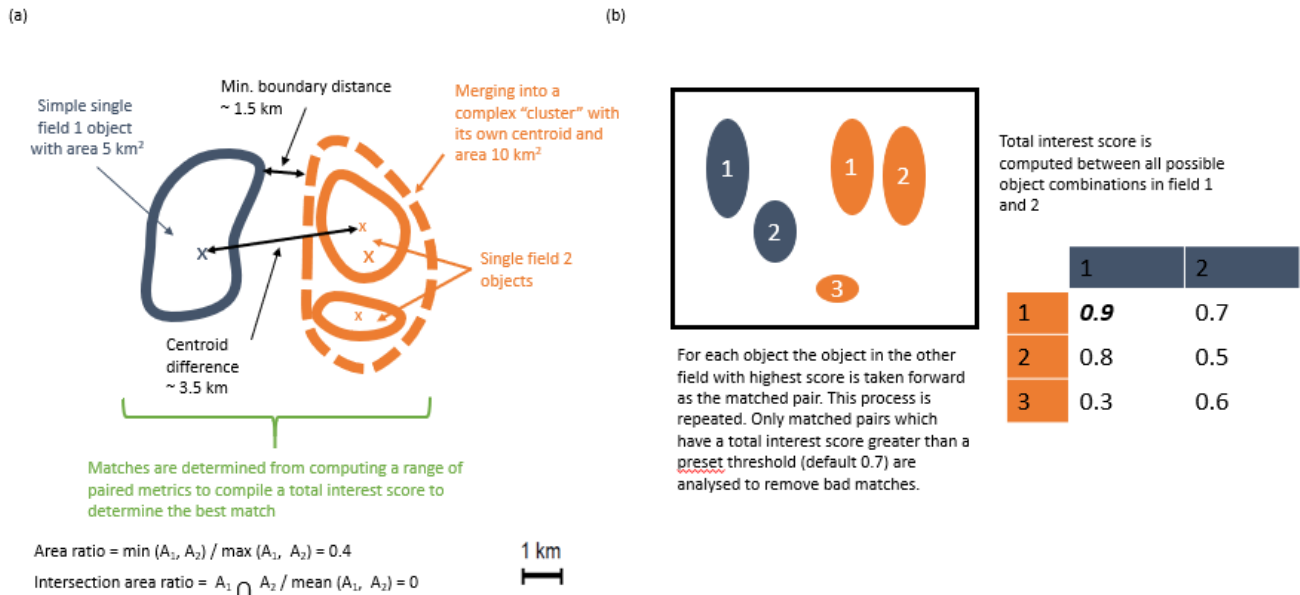
190 This section provides a brief description of the Method for Object-Based Diagnostic Evaluation  
191 (MODE), first described in Davis et al. (2006) and its extension MODE Time-Domain (MTD).  
192

193 MODE and MTD can be used on any temporal sequence of two gridded data sets which contain features  
194 that are of interest to a user (whoever that user may be, model developer or more applied). By extracting  
195 only the feature(s) of interest, the method allows one to mimic what humans do, but in an objective  
196 way. Once identified the features can then be mathematically analysed over many days or seasons to  
197 compute aggregate statistics of behaviour. MODE can be used in a very generalised way. The key  
198 requirements are to 1) have gridded fields to compare and 2) be able to set a threshold for identifying  
199 features of interest.

200

201 In this instance the comparison will involve the AMM7v11 model data assimilation analysis and the  
202 gridded L4 satellite product. MODE identifies the features (called objects), as areas for which a  
203 specified threshold is exceeded, here it is a Chl-*a* concentration. Consider Figure 2 which shows a  
204 number of objects that have been identified after a threshold has been applied to two fields (blue and  
205 orange). The identified objects in the two fields are of different sizes and shapes and do not overlap in  
206 space, though they are not far apart. Object characteristics or attributes such as the area and mass-  
207 weighted centroid are computed for each single object. Simple (also known as single) objects can be  
208 *merged* (to form clusters) within *one* field (illustrated here for the orange field). This may be useful to  
209 do if it is clear that there are many small objects close together which should really be treated as one.  
210 Furthermore, objects in one field can be *matched* to objects in the other field. To find the best match an  
211 interest score is computed for each possible pairing. The components used for computing the interest  
212 score can be tuned to meet specific user needs. In (a) it is based on the area ratio, intersection-area ratio,  
213 minimum boundary distance and centroid difference. Furthermore, the components can be weighted  
214 according to relative importance. Given a scenario where there are 2 identified objects in the blue field  
215 and 3 in the orange field (b) shows the interest score for each possible pairing in this hypothetical  
216 example. Only the pairing with the highest score is analysed further, and only if it exceeds the set  
217 threshold for defining an acceptable match. The default value for this is 0.7. Once these matches are

218 completed summary statistics describing the objects (both matched and unmatched) and matched object  
 219 pairs are produced. These statistics can be used to identify similarities and differences between the  
 220 objects identified in two different data sets, which can provide diagnostic insights on the relative  
 221 strengths and weaknesses of one compared to the other.



222  
 223 **Figure 2 Schematic illustrating some of the key components of identifying objects using MODE. (a) Defining some of**  
 224 **the terminology and key components for computing matched pairs. (b) Example of how the best matched pair is**  
 225 **identified.**  
 226

227 The important steps for applying MODE can be summarised as follows (which are described in detail in  
 228 Davis et al. 2006):

- 229 1) Both forecast and observation (or analysis) need to be on the same grid. Typically, this means  
 230 interpolating the observations to the model grid to avoid the model being expected to resolve  
 231 features which are sub-grid scale.
- 232 2) Depending on how noisy the fields are they should be smoothed. It is worth remembering that  
 233 the numerical discretisation implies that any model's true resolution (i.e. the scales which the  
 234 model is resolving) is between 2 and 4 times the horizontal grid (mesh) resolution. The number  
 235 of objects identified will vary inversely with the smoothing radius.

- 236 3) Define a threshold which captures the feature of interest and apply it to both the smoothed  
237 forecast and observed fields to identify simple objects as shown in Figure 2.
- 238 4) Any smoothing is only for object identification purposes. The original intensity information  
239 within the object boundaries is analysed.
- 240 5) Lastly, the object matching is accomplished using a fuzzy logic engine (low level artificial  
241 intelligence), which is expressed as the so-called “interest” score as shown in Figure 2(b). The  
242 higher the score the stronger the match. All objects are compared in both fields and interest  
243 scores are computed for all combinations. A threshold is set on the interest score value (typically  
244 0.7) to denote which are the best matches, and on the unique pairing with the highest score is  
245 kept for analysis purposes. Some objects will remain unmatched (either because there is none or  
246 because there are no interest values above the set threshold to provide a credible match) and  
247 these can be analysed separately.

248 MODE is highly configurable. To gain an optimal combination of configurable parameters for each  
249 application requires extensive sensitivity testing to gain sufficient understanding of the behaviour of the  
250 data sets to be examined, and to achieve, on average, heuristically the right outcome. Initial tuning  
251 requires user input to check whether the method is replicating what a human would do.

- 252 1) The sensitivity to threshold and smoothing radius should be explored. The threshold and  
253 variability in the fields can affect the number of objects which are identified. The process of  
254 exploring the relationship between threshold and smoothness helps to identify what would  
255 heuristically be considered a reasonable number of objects.
- 256 2) The sensitivity to the merging option must also be investigated. In this instance the merging  
257 option had very little impact.
- 258 3) The behaviour of the matching can also be configured, with a number of options ranging from  
259 the simple to the more complicated, which added computational expense. There may be very  
260 little difference in outcomes, but it is worth checking. Here the *merge\_both* option was used but  
261 it was not strictly necessary as there was little difference between the available options.

262

263 Note also that a minimum size (area) is set for object identification. This is often a somewhat pragmatic  
264 choice. If the size is set too small, too many objects are identified, which end up being merged. If too  
265 large, very few objects are identified. Here a minimum area of 10 grid squares ( $\sim 70 \text{ km}^2$ ) was used for  
266 an object to be included in the analysis. For this study the default settings were used for matching and  
267 computing the interest score (as provided in the default configuration file (see example configuration  
268 files in [https://github.com/dtcenter/MET/tree/main\\_v8.1/met/scripts/config](https://github.com/dtcenter/MET/tree/main_v8.1/met/scripts/config)). The default threshold of  
269 0.7 for the interest score was also used to identify acceptable matches.

270

271 Identical to MODE, identifying time-space objects in MTD uses smoothing and thresholding. Applying  
272 a threshold yields a binary field where grid points exceeding the defined threshold are set to one. At this  
273 stage each region of non-zero grid points in space and time is considered a separate object, and the grid  
274 points within each object are assigned a unique object identifier. For MTD the search for contiguous  
275 grid points not only means examining adjacent grid points in space, but also the grid points in the same  
276 or similar location at adjacent times to define a space-time object. The same fuzzy logic-based  
277 algorithms used for merging and matching in MODE apply to MTD as well. Similarly, to MODE a  
278 minimum volume must be set. Here a volume threshold of 1000 grid squares was imposed for space-  
279 time object identification to be included in the analysis. This represents the accumulated number of grid  
280 squares associated with an object over consecutive time slices. Otherwise, the default settings were used  
281 for object matching. For MTD a lower interest score of 0.5 was used for matching objects. Finally, it is  
282 worth noting that the MODE and MTD tools, though similar, are completely independent of each other,  
283 and were set up differently here. MODE is ideal for understanding the identified features in individual  
284 daily fields in some detail. MTD, it was felt, would be best used to look at larger scales. Here it was set  
285 up to capture the most significant (in size) and long-lasting blooms.

286

### 287 **3.2 Defining Chl-*a* concentration thresholds and other choices on tuneable parameters**

288 Chl-*a* can vary over several orders of magnitude. Often  $\log_{10}$  thresholds are used to match the fact that  
289 Chl-*a* follows a lognormal distribution (e.g. Campbell 1995). Defining thresholds can be difficult: on  
290 the one hand there is the desire to capture events of interest, so the thresholds should not be too low,

291 whereas on the other hand if the thresholds are too high no events are captured and there is nothing to  
292 analyse. From a regional perspective the values of interest are typically in the range of 3–5 mg m<sup>-3</sup>  
293 (Schalles, 2006), though higher values are present. For this study, to set of equally spaced logarithmic  
294 thresholds, ranging between 0.2 and 1.4 log<sub>10</sub>mg m<sup>-3</sup> were applied to the Chl-*a* fields, corresponding to  
295 Chl-*a* concentrations between 1.62 and 25 mg m<sup>-3</sup>. Doing this removed the need to transform the data.  
296 In the paper the primary focus is on the 2.5 mg m<sup>-3</sup> threshold, though some results for the 4 and 6.3 mg  
297 m<sup>-3</sup> are also presented.

298

299 In addition to the interpolation of the L4 ocean colour product onto the AMM7 grid, it is important to  
300 ensure that MODE and MTD use optimal settings for the fields under study. Results are sensitive to  
301 characteristics of the fields (how smooth or noisy). Right at the start the emphasis was on finding the  
302 right combination of Chl-*a* concentration threshold and smoothing, balancing the need for identifying  
303 objects with keeping the number of objects manageable. The guiding principles in identifying the right  
304 combination were to ensure that the daily object count remained less than 30. Furthermore, the  
305 smoothing applied needs to be reduced with increasing concentration thresholds because objects  
306 become smaller and are less frequent. This is to ensure that too much smoothing does not remove more  
307 intense objects from the analysis. However, pushing the concentration threshold too high may also be  
308 too detrimental; identified objects may be spurious and too few objects will mean meaningful statistics  
309 cannot be compiled. AMM7v11 analyses are on a ~7 km grid.

310

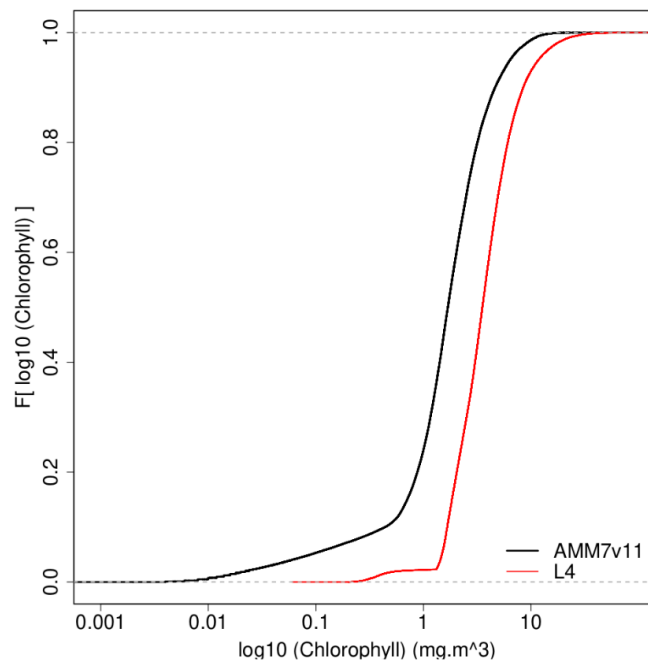
311 For the lowest thresholds including 2.5 and 4.0 mg m<sup>-3</sup> a smoothing radius of 5 grid squares (~35 km)  
312 was applied to both L4 and AMM7v11 fields, but for higher thresholds (e.g. 6.3 mg m<sup>-3</sup>) the smoothing  
313 radius was reduced to 3 grid squares, to prevent the higher peak concentrations, which are often small in  
314 spatial extent, from being lost due to the smoothing. Thresholds above 6.3 mg m<sup>-3</sup> yielded too few  
315 objects to be analysed with any rigour. The smoothing was particularly necessary for the L4 product  
316 which, because of its native 1 km resolution is able to resolve very small (noisy) objects typically found  
317 near the coast and which a 7 km resolution model cannot resolve. For the MTD analysis, objects in the

318 L4 ocean colour product and the AMM7v11 analyses were defined using a Chl-*a* concentration  
319 threshold of 2.5 mg m<sup>-3</sup>.

## 320 4. Results

### 321 4.1 Chl-*a* distributions

322 It is important to understand the nature of the underlying L4 and AMM7v11 Chl-*a* distributions and any  
323 differences between them. This can be done by creating cumulative distribution functions (CDF) of the  
324 log<sub>10</sub> L4 and AMM7v11 Chl-*a* concentrations, by taking all grid points in the domain and all dates in  
325 the study period. These are plotted in Figure 3, showing that there is an offset between the distributions,  
326 the AMM7v11 analysis having more low concentrations, though the distributions appear to be  
327 converging in the upper tail.

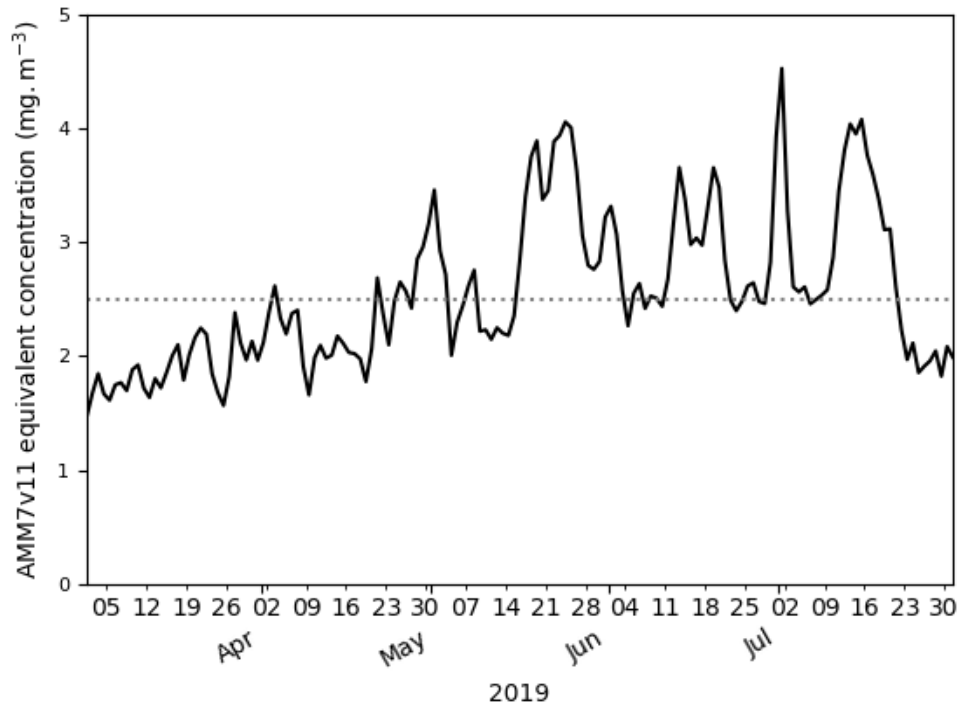


328

329 **Figure 3 Empirical cumulative distribution functions of the log<sub>10</sub> Chl-*a* concentration for the L4 ocean colour**  
330 **product and AMM7v11 analyses for the 2019 bloom season.**

331 Exploring this further the AMM7v11 and L4 Chl-*a* concentration CDFs can be derived for each  
332 individual day, rather than for the season as a whole. From these the centile where the L4 product is less

333 than equal to  $2.5 \text{ mg m}^{-3}$  can be compared to the corresponding AMM7v11 centile value. The daily  
334 matched centile Chl-*a* values provide an estimate of the daily bias. This is plotted in Figure 4 as a time  
335 series for the 2019 bloom season. It shows that the daily AMM7v11 corresponding centile values are  
336 mainly in the range of  $\sim 1.5\text{--}4.5 \text{ mg m}^{-3}$ , averaging out to  $2.9 \text{ mg m}^{-3}$  over the season, which suggests a  
337 modest difference overall. The larger day-to-day variations show some cyclical patterns. There are  
338 notable peaks at the end of May and the beginning of July. An inspection of the fields (not shown)  
339 suggests that at these times the AMM7v11 appears to have higher Chl-*a* concentrations over large  
340 portions of the domain compared to the L4 product.



341  
342 **Figure 4** The day-to-day AMM7v11 centile Chl-*a* value corresponding to the L4 product centile representing  $2.5 \text{ mg}$   
343  $\text{m}^{-3}$  derived from the L4 daily CDFs. The mean AMM7v11 Chl-*a* equivalent centile value for the 2019 season is  $2.9 \text{ mg}$   
344  $\text{m}^{-3}$ .

345  
346 In employing a threshold-based approach, generally the same threshold is applied to both data sets. In  
347 the presence of a bias this requires a little bit of thought. In extreme cases, it could mean the inability to  
348 identify objects in one of the data sets, which would then mean objects cannot be matched and paired,  
349 negating the purpose of a spatial method like MODE or MTD. Not being able to identify any objects

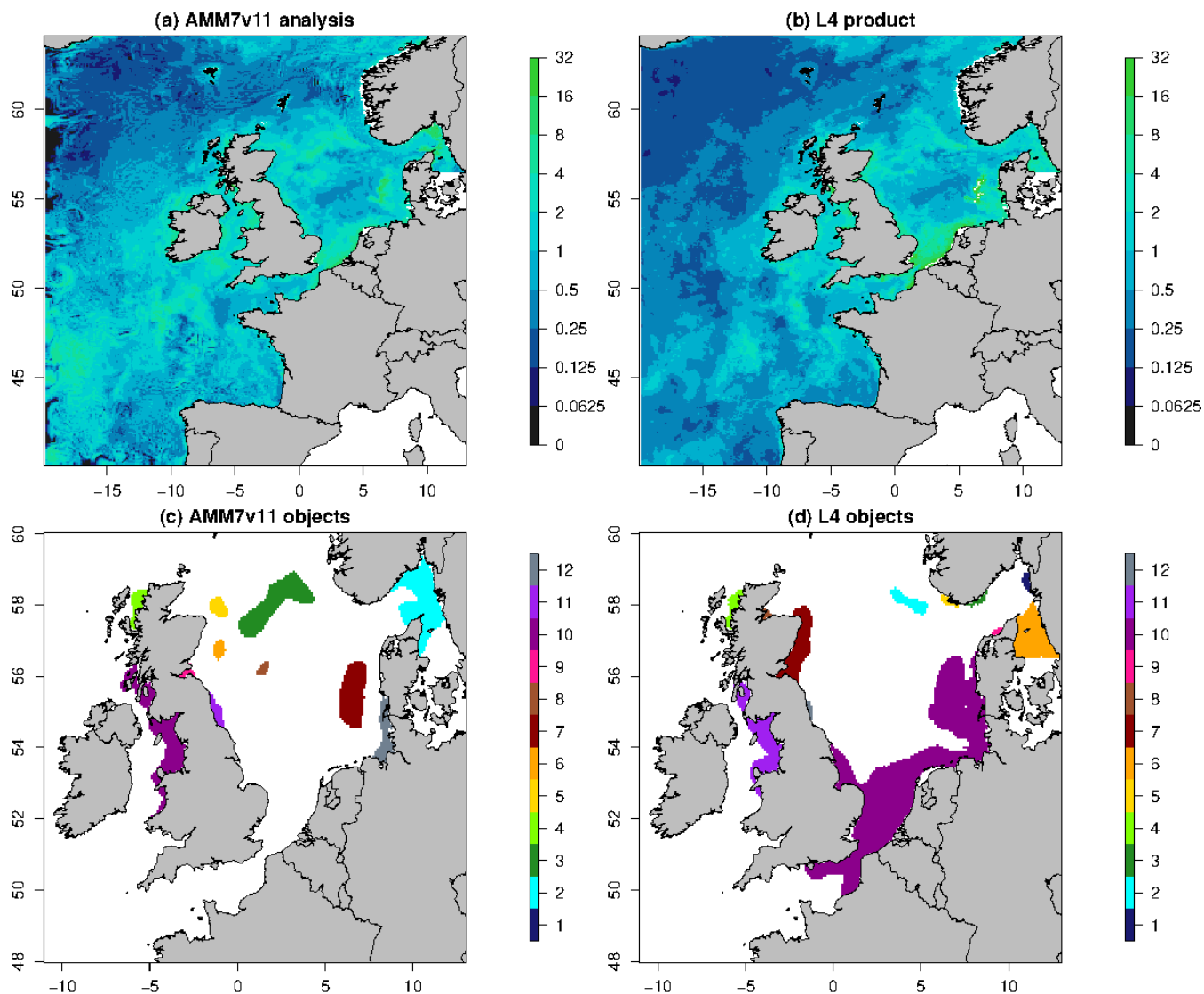
350 does provide some useful information, though arguably not enough context. The lack of objects does  
351 suggest the presence of a bias but it does not provide any sense of whether the model is producing a  
352 constant value of Chl-*a* for example, which would be of no use to the user, or whether it does capture  
353 regions of enhanced Chl-*a*, albeit with an offset which means it does not exceed the set threshold.  
354 Therefore, a more likely scenario is that a bias could partially mask relevant signals in the derived  
355 object properties, which could lead to the potential misinterpretation of results. If there is a significant  
356 risk of this occurring the bias could be addressed before features are identified to ensure that the  
357 primary purpose of using a feature-based assessment can be achieved, i.e. identifying features of interest  
358 in two sets of fields to assess their location, timing and other properties and assessing their skill. The  
359 fact that there is an intensity offset should not prevent the method from providing information about the  
360 skill of identified features. In this instance, though there is bias, it did not prevent the identification of  
361 objects in either fields to the extent where the results did not reflect the potential for the analyses to  
362 provide features which could be matched, paired and compared.

## 363 **4.2 Visualising daily objects**

364 Figure 5 shows the daily Chl-*a* concentration fields as represented in the L4 ocean colour product and  
365 the AMM7v11 analyses for 21 April 2019, which is near the peak of the bloom season. The respective  
366 fields are plotted in (a) and (b), noting that the 1 km resolution L4 product has been interpolated onto  
367 the ~7 km AMM7 grid. Applying a threshold of  $6.3 \text{ mg m}^{-3}$  to both with a smoothing radius of ~21 km  
368 (3 grid lengths) yields 8 objects in the AMM7v11 analysis (7 visible in this zoomed region) and 11  
369 objects in the L4 product. As discussed, the bias described in Section 4.1 does not appear to prevent the  
370 identification of objects in the L4 product and the AMM7v11 analyses, and the process of finding  
371 matches is possible.

372





373

374 **Figure 5 Daily Chl-*a* concentrations (in  $\text{mg m}^{-3}$ ) for 21 April 2019: (a) AMM7v11 analysis and (b) L4 ocean colour**  
 375 **product. The MODE objects shown in (c) and (d) are identified using a threshold of  $6.3 \text{ mg m}^{-3}$  and a smoothing**  
 376 **radius of  $\sim 21 \text{ km}$ . The colour matches the object identification number.**

377

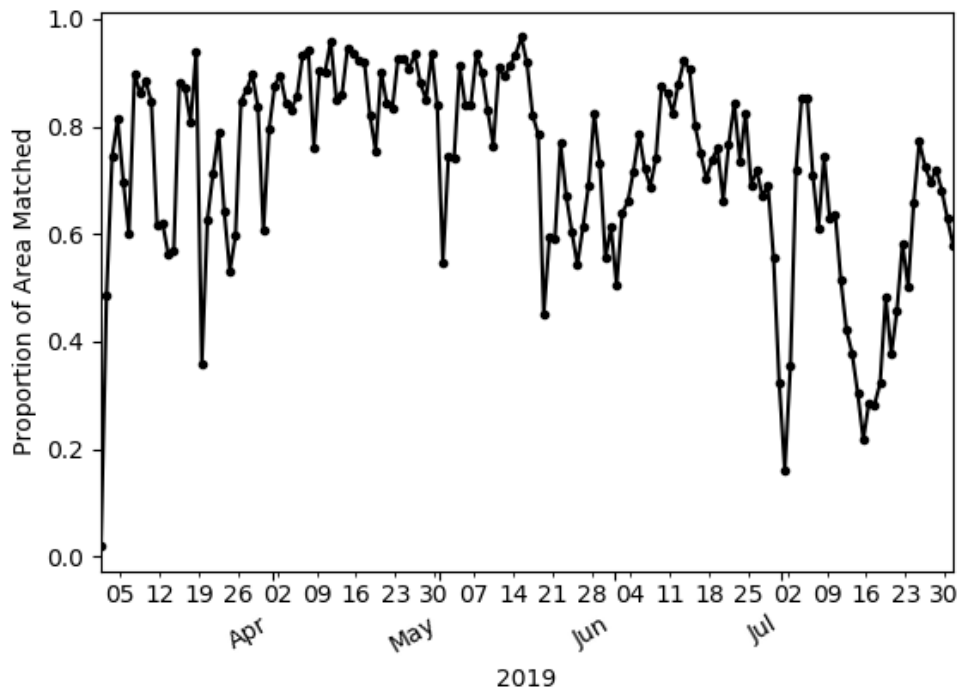
### 378 4.3 Spatial characteristics

379 This section demonstrates the kinds of results that can be extracted from the two-dimensional MODE  
 380 objects. Aspects of the marginal (AMM7v11 or L4 product only) and joint (matched/paired)

381 distributions can be examined. This includes object size (as a proxy for area) but also the proportion of  
382 areas that are matched or unmatched.

383

384 Firstly, how similar is the L4 ocean colour product and the AMM7v11 analysis in terms of the features  
385 of most interest, i.e. the Chl-*a* blooms? Figure 6 shows the evolution of the proportion of matched  
386 object areas (to total combined area) through the 2019 season, when using MODE to compare the L4  
387 product and AMM7v11 analyses, to further explore the differences (and similarities) between them. A  
388 value of one would suggest that all identified areas are matched. Values less than one suggest that some  
389 objects remain unmatched. The relatively high values of matched object-to-total area during April are  
390 due to the large numbers of well-matched, physically small coastal objects in addition to the larger Chl-  
391 *a* bloom originating in the Dover Straits (not shown). There is a notable minimum at the beginning of  
392 July. Inspecting the MODE graphical output reveals this is in part due to only a few small objects being  
393 identified, and this is compounded by their complete mismatch; the L4 objects are all coastal, whilst the  
394 AMM7v11 objects are either coastal (but not in the same location as L4 objects) or in the deep waters of  
395 the North Atlantic, to the north-west of Scotland. The relatively high proportions either side of this time  
396 arise from a better correspondence in placement of the coastal objects (noting that there is a distance  
397 limit on how far objects can be apart for the matching process to have a positive contribution to the  
398 interest score).



399

400 **Figure 6 Proportion of total object area which is matched. Underlying matched and unmatched object areas (in units**  
 401 **of numbers of grid squares) are taken from the MODE output. These areas are for the 2.5 mg m<sup>-3</sup> concentration**  
 402 **threshold objects.**

403

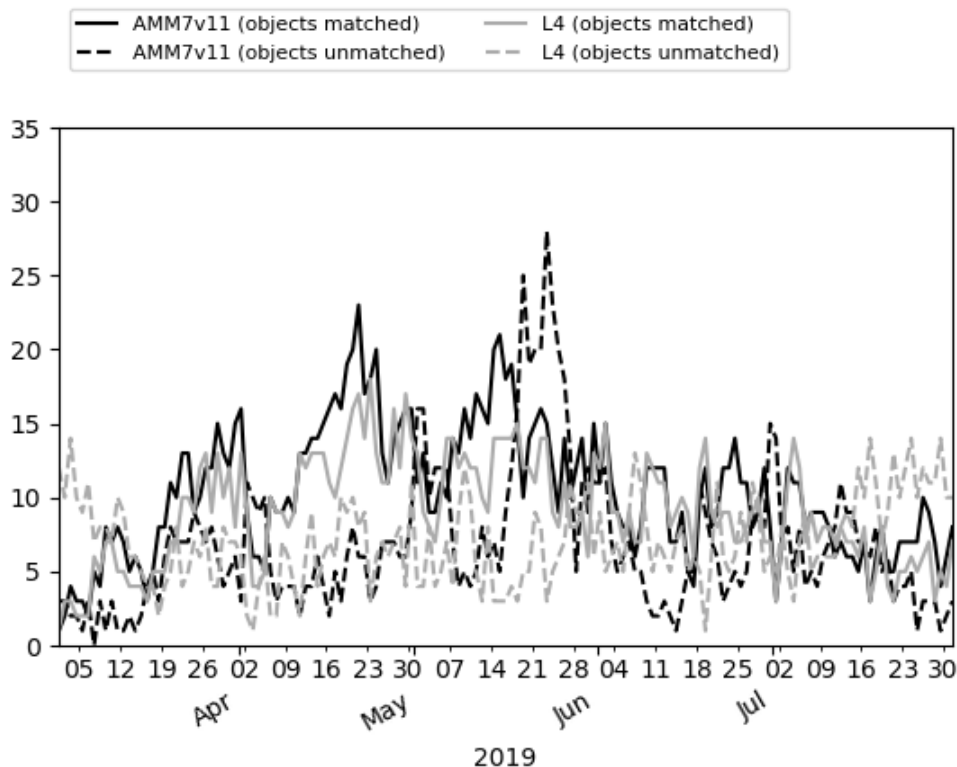
404 Overall, the AMM7v11 analysis is similar, but clearly not identical, to the L4 product. Best  
 405 correspondence appears to be during the first half of the bloom season. Later in the season the model's  
 406 determination to produce blooms in deep North Atlantic waters is a model deficiency that the  
 407 assimilation is (at this stage) unable to fix. The AMM7v11 analyses could conceivably be used as a  
 408 credible source for assessing the AMM7 Chl-*a* forecasts in the future. The major benefit of using a  
 409 model analysis is that it is at the same spatial resolution, with the same ability to resolve Chl-*a* bloom  
 410 objects, especially along the coast (i.e. the analysis limits the uncertainty due to whether an object could  
 411 be missing due to the inability of the model to resolve the feature).

412

413 The day-to-day number of objects identified through the 2019 bloom season is shown in Figure 7,  
 414 illustrating how elements of the marginal and joint distribution information provided by MODE can be

415 used together. Here both matched (joint) and unmatched (marginal) objects are shown. From an  
416 interpretation perspective there should be fewer unmatched objects than matched ones (ideally there  
417 would be no unmatched objects in either the forecast or the analysis). In Figure 7 the number of objects  
418 in AMM7v11 starts off small and increases as the bloom develops. For the L4 product there are already  
419 many objects identified at the start of the timeseries, leading to many unmatched L4 objects. A spike in  
420 the number of matched objects seen in early April can be attributed to several coastal locations, which  
421 appear to be spatially well-matched. In addition, a larger Chl-*a* bloom is seen in the Dover Straits region  
422 in the L4 product and although not exactly spatially collocated, the objects are matched. There are a  
423 consistently large number of unmatched objects seen in the AMM7v11 analysis and L4 ocean colour  
424 product from the end of May onwards. In the AMM7v11 analysis this appears to be due to an increase  
425 in small objects identified, mainly to the west, north and east of the United Kingdom. The increase in  
426 unmatched objects in the L4 ocean colour product is of a different origin, being due to an increase in  
427 localised coastal blooms. Generally, the AMM7v11 analyses do not have the resolution to resolve these.  
428 Overall, there are 2632 AMM7v11 bloom objects identified in the season using the  $2.5 \text{ mg m}^{-3}$   
429 threshold, and 2341 L4 bloom objects, with 56% of AMM7v11 objects matched and 59% of L4 objects  
430 matched.

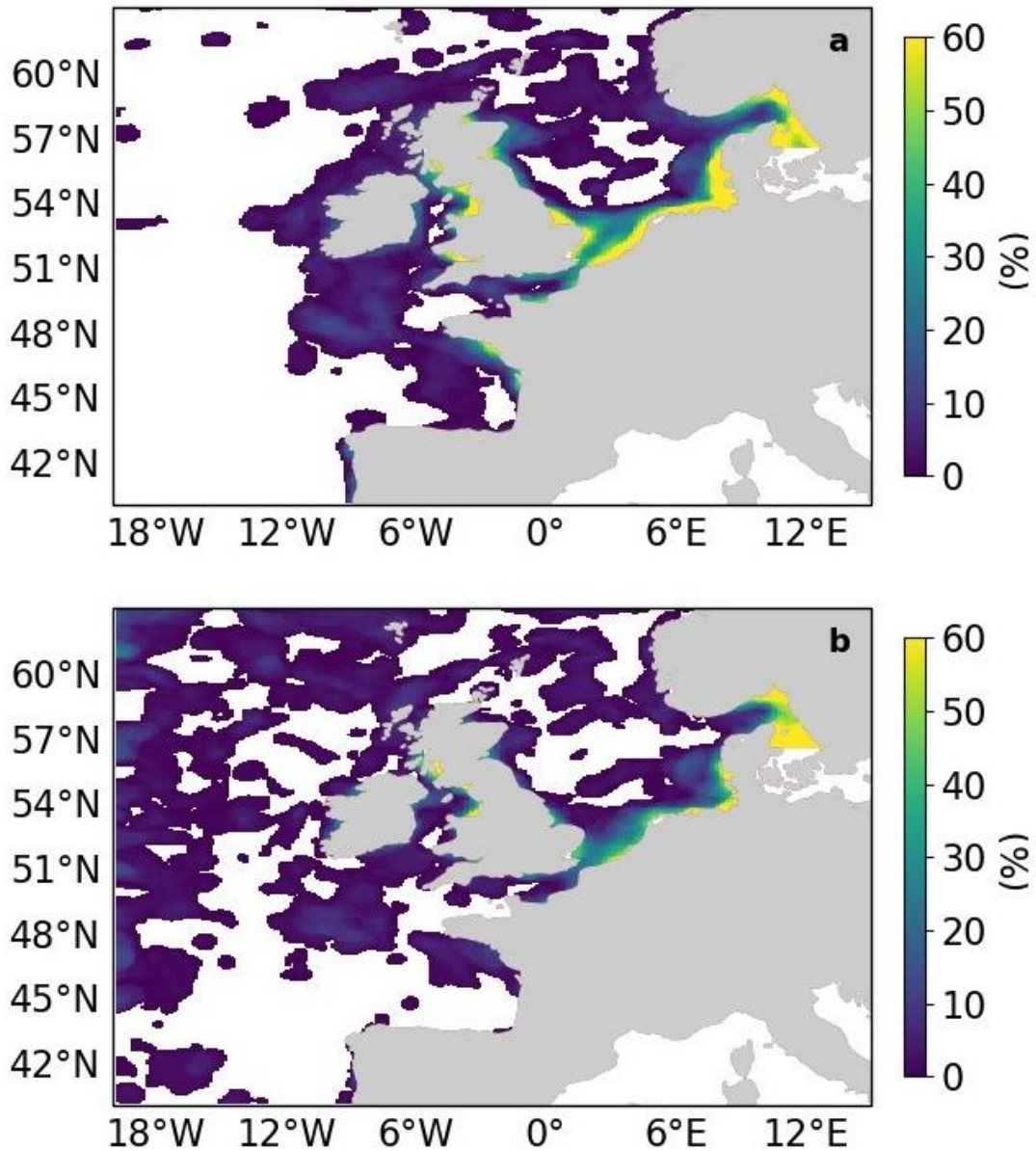
431 The identified objects in AMM7v11 and the L4 product can also be considered spatially over the season  
432 by compositing the objects. This is done by counting the frequency with which a given grid square falls  
433 within an identified object on any given day, essentially creating a binary map. These can be added up  
434 over the entire season to produce a spatial composite object or temporal “frequency-of-occurrence” plot.



435

436 **Figure 7 Time series of the number of matched and unmatched objects from MODE comparing AMM7v11 analyses**  
 437 **(black) with L4 satellite product (grey). Objects are identified using a threshold of  $2.5 \text{ mg m}^{-3}$ . Total object numbers**  
 438 **for the season are 2341 for L4 satellite product and 2632 for AMM7v11.**

439 Figure 8 shows this spatial composite for the 2019 bloom season for the L4 ocean colour product  
 440 objects (a) and the AMM7v11 objects (b). These are the composites based on the  $2.5 \text{ mg m}^{-3}$  threshold  
 441 objects. There are areas, for example in the South West Approaches, where there appears to be a good  
 442 level of consistency. AMM7v11 analyses have elevated Chl-*a* values along the northern and western  
 443 edges of the domain, for a low proportion of the time, which are not seen in the L4 product. This is  
 444 likely due to the way that nutrient and phytoplankton boundary conditions are specified in AMM7v11.  
 445 Overall, the low temporal frequency extent of the AMM7v11 objects is greater than for the L4 product.

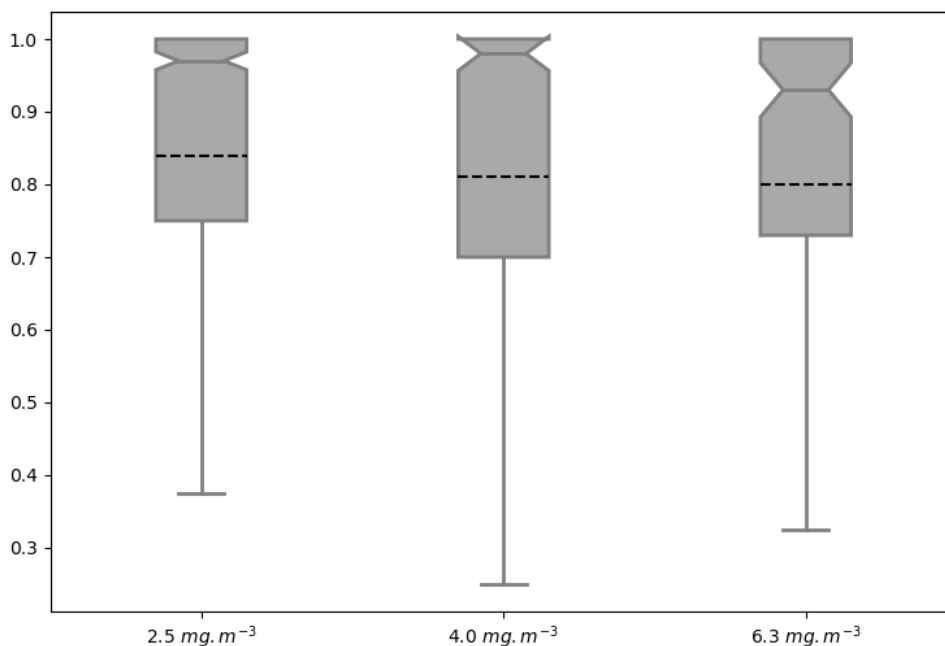


446

447 **Figure 8 Object composites (the proportion of time for which an object was present at the grid box throughout the**  
 448 **2019 bloom season) for (a) the L4 ocean colour product objects and (b) the AMM7v11 analysis objects.**

449 Thus far all the attributes have been based on only the AMM7v11 or L4 objects. The distribution of  
 450 object properties, derived for the season from the daily comparisons, can be summarised using box-and-  
 451 whisker plots. Recall that the box encompasses the inter-quartile range (IQR, 25<sup>th</sup> to 75<sup>th</sup> percentile) and

452 the notch and line through the box denotes the median or 50<sup>th</sup> percentile. The dashed line represents the  
453 mean, and the whiskers show  $\pm 1.5$  times the IQR. For clarity, values outside that range have been  
454 filtered out of the plots shown here. Figure 9 shows the intersection-over-area paired object attribute  
455 distribution as box-and-whisker plots for all object pairs during the 2019 bloom season, comparing the  
456 AMM7v11 analyses to L4 for three of the thresholds: 2.5 and 4.0 and 6.3 mg m<sup>-3</sup>. The intersection-over-  
457 area diagnostic gives a measure of how much the matched (paired) objects overlap in space. If the  
458 objects do not intersect, this metric is 0. The ratio is bounded at 1 because any area of overlap is always  
459 divided by the larger of the two object areas. The IQR for the 2.5 mg m<sup>-3</sup> threshold is 0.25 with 50% of  
460 paired objects having an intersection-over-area of 0.97 or greater. However, the lower whisker spans a  
461 large range of values to as low as 0.375, suggesting that there is a proportion of object pairs with only  
462 small overlaps. There is quite a difference between the median (notch) and the mean (dashed line) for  
463 this metric, suggesting the distribution is skewed with the mean affected more by many small overlaps.  
464 For the 4.0 mg m<sup>-3</sup> threshold paired objects the intersection-over-area distribution is much broader,  
465 though the difference between the mean and medians is similar. The proportion of paired objects with  
466 smaller overlaps has also increased. This should not be surprising given that the objects generally get  
467 smaller with increasing threshold such that the ability for object pairs to overlap actually decreases  
468 unless they are very closely collocated. At the 6.3 mg m<sup>-3</sup> threshold the median is lower (0.93) with a  
469 similar difference from the mean, however the sample size is much smaller (only 130 paired objects  
470 over the season).



471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

**Figure 9** Box-and-whisker plots of the paired object property “intersection area” ratio computed by dividing the spatially collocated area between the paired objects by the largest of either the AMM7v11 or L4 observed object areas (to keep the ratio to be bounded by 0 and 1). Three object thresholds are shown: 2.5  $\text{mg.m}^{-3}$ , 4.0  $\text{mg.m}^{-3}$  and 6.3  $\text{mg.m}^{-3}$ . Smoothing radii of 5, 5 and 3 grid lengths were applied for the three thresholds respectively. The sample sizes for each threshold were 1004, 401 and 130 paired objects respectively.

#### 4.4 Incorporating the time dimension

Having information in space *and* time enables one to ask, and hopefully answer questions such as: “*did the model predict the bloom to start in the observed location?*” or “*did the model predict the onset at the right time?*” and “*did the model predict the peak (in terms of extent) and duration of the bloom correctly?*”.

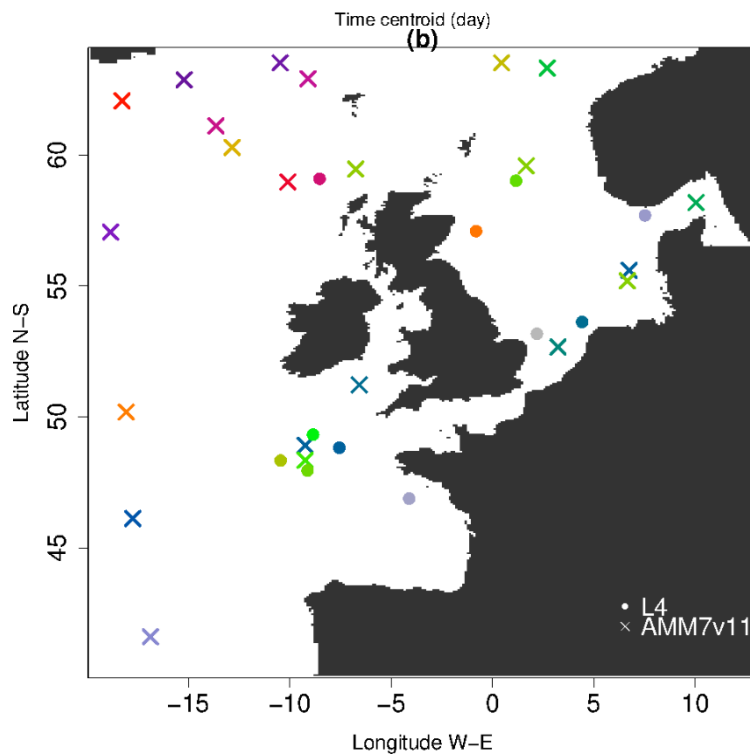
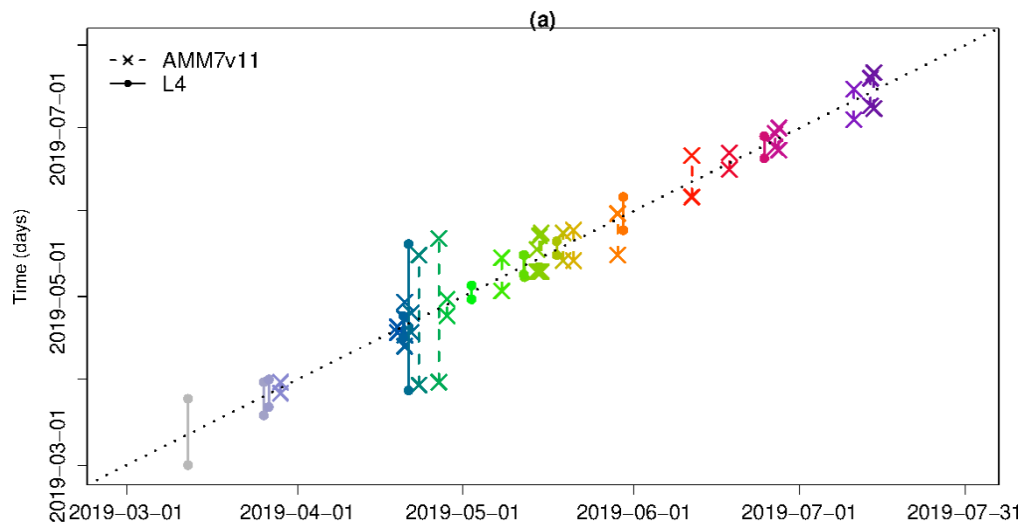
MTD identifies objects in space and time. As previously described, all MTD results are based on a 2.5  $\text{mg.m}^{-3}$  threshold applied to both the L4 ocean colour products and AMM7v11 analyses. A time centroid is derived from a time series of the spatial (two-dimensional) centroids which are computed for each (daily) time slice. In addition to this, each identified MTD object has a start and end time, and a



489 geographical location of the time centroid, which is the average of the two-dimensional locations. The  
490 time component of the time centroid is weighted by volume.

491

492 The temporal progression of the 2019 bloom season as defined by the MTD objects' start and end times  
493 as well as the date of their time centroids is shown in Figure 10, providing a clear view of the onset and  
494 demise of each object (bloom episode). In total there are 22 AMM7v11 and 11 L4 MTD objects. The x-  
495 axis in (a) represents elapsed time. The location of the vertical lines along the x-axis on any given date  
496 indicates the date of the time centroid whilst the duration of the space-time object can be gleaned from  
497 the y-axis based on the start and end of the vertical line which defines the time the object was in  
498 existence. Solid lines represent the L4 product objects whereas dashed lines represent the AMM7v11  
499 objects. The colour palette is graduated from grey and blue through green, yellow, red, and purple,  
500 denoting the relative time in the season. In (a) the first Chl-*a* bloom object in the AMM7v11 analysis  
501 was identified on 29 March 2019 whereas in the L4 ocean colour product this was on 3 March, 26 days  
502 earlier. The first time the L4 product and AMM7v11 analyses have concurrent objects (blooms) is in  
503 late March. The L4 product also suggests that the season ends 30 June whereas the AMM7v11 analyses  
504 persists the bloom season with objects identified until 23 July. Most AMM7v11 objects are of relatively  
505 short duration, but overall, most groups of AMM7v11 objects have some temporal association with an  
506 L4 product object around the same time, though this does not mean they are geographically close to  
507 each other. This is illustrated in Figure 10(b) which provides the spatial context to (a). The colours and  
508 symbols are consistent for (a) and (b) and show that even when the MTD objects are identified at the  
509 same time they may be geographically quite far apart, or more typically there is no L4 counterpart  
510 (filled circle) to an AMM7v11 bloom object (cross). The north- and westward progression of the bloom  
511 as the season unfolds can be seen through the use of the colours, with the AMM7v11 analysis producing  
512 many more objects in deeper waters to the north and west of the domain.



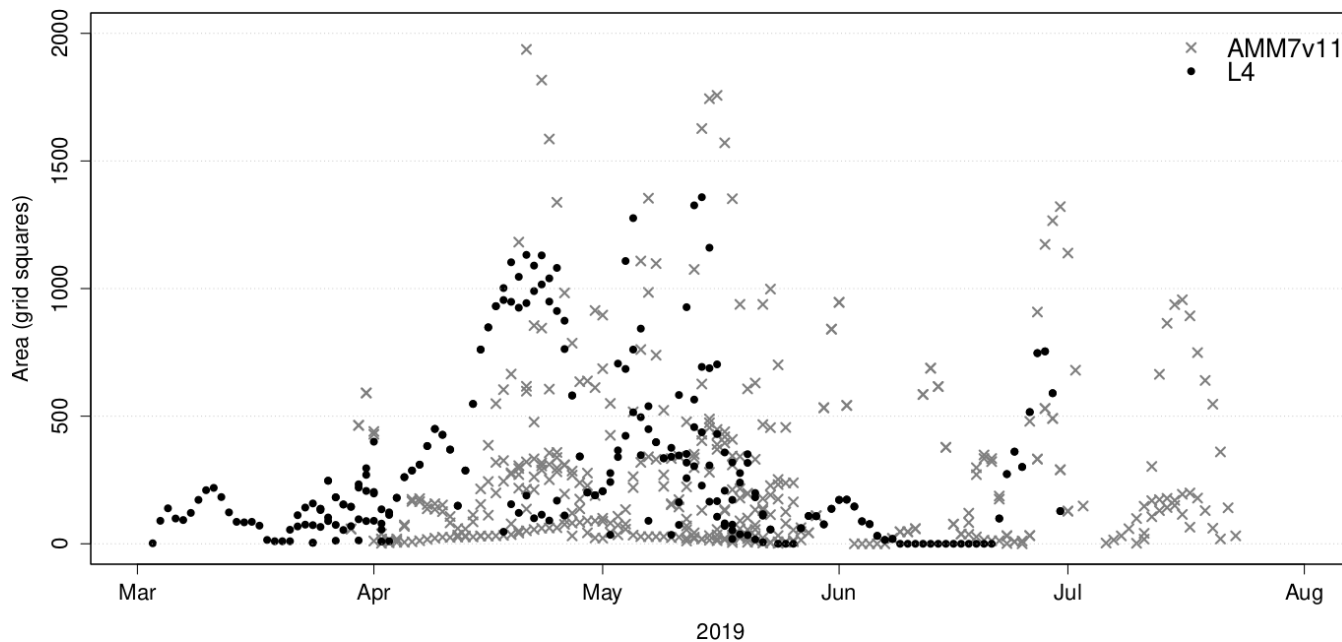
513

514

515 **Figure 10** Space-time information from the L4 (filled circle) and AMM7v11 (cross) MTD objects. (a) The timing of  
 516 each identified bloom event (time centroid) plotted on the x-axis against the duration of the bloom event, denoted by  
 517 the vertical line which represents the start and end time of each space-time object. The colours provide the ability to  
 518 track the relative location within the 2019 season. (b) Spatial location of the time centroid shown in (a) to indicate  
 519 that even if AMM7v11 and L4 objects exist at the same time they may not be geographically close. Colours are  
 520 coordinated between (a) and (b).

521

522 In this instance it is also illuminating to consider a time series of all identified daily object areas (which  
523 are used to compute the volume of MTD objects). These are plotted in Figure 11 showing all daily L4  
524 object areas in black filled circles, and the AMM7v11 object areas in grey crosses. The main purpose is  
525 to highlight the relative size of the L4 and AMM7v11 objects on any given day, as well as how many  
526 objects there were. Recall that these are the objects identified using a Chl-*a* concentration threshold of  
527  $2.5 \text{ mg m}^{-3}$ . Some of the AMM7v11 objects are considerably larger than those in L4 though in the  
528 middle part of the bloom season between mid-May and end June there is reasonable correspondence in  
529 identifying the peak in terms of extent and activity, just not necessarily at exactly the same time or  
530 location. Of course, the AMM7v11 areas may also be larger because of the difference in the  
531 distributions noted in Figure 3, one of the reasons an awareness of the presence of any biases is  
532 important when interpreting results. As seen in Figure 10(b), the area time series also illustrates the  
533 offsets in the start and end of the bloom season. Some of the objects detected in AMM7v11 beyond the  
534 end of the observed bloom season provided by L4, suggests that at least three substantial areas are still  
535 diagnosed to exceed the threshold of  $2.5 \text{ mg m}^{-3}$  into July. Taking the start of the earliest space-time  
536 object as the onset of the bloom season and the end of the last object as the end, the 2019 season is 119  
537 days long based on the L4 product, and 117 days in the AMM7v11 analysis. Therefore, the overall  
538 length of the season as defined by the space-time objects is comparable in the AMM7v11 analysis,  
539 albeit with a substantial offset.



540

541 **Figure 11 Time series of all identified single simple MTD object areas in the AMM7v11 analysis and the L4 ocean**  
 542 **colour product.**

543 With only 22 AMM7v11 and 11 L4 product MTD objects, which are temporally and geographically  
 544 well dispersed, three of the L4 objects remained unmatched, leaving only 8 matched MTD objects for  
 545 the 2019 bloom season with an overall interest score greater than 0.5. This represented an insufficient  
 546 sample for drawing any robust statistical conclusions. Nevertheless, some inspection of the paired MTD  
 547 object attributes are summarised below:

- 548 • The spatial centroid (centre of mass) differences can be extensive, but the majority are within 0 to  
 549 100 grid squares apart (i.e. up to ~700 km).
- 550 • The majority of paired objects have time centroid differences +/- 10 days.
- 551 • Considering the volumes of the space-time objects, half the paired objects have volume ratios of less  
 552 than 1, i.e. AMM7v11 objects tend to be smaller or similar in size. The other pairs have ratios as  
 553 high as 4.
- 554 • Overlaps between AMM7v11 and L4 MTD objects remain small and infrequent with only one pair  
 555 with a significant overlap in space and time.

## 556 5. Discussion and conclusions

557 MODE and MTD were used as two distinct but related feature-based diagnostic verification methods to  
558 evaluate and compare the pre-operational AMM7v11 European North West Shelf Chl-*a* concentration  
559 bloom objects to those identified in the satellite-based L4 ocean colour product. Nominally blooms were  
560 said to occur when the concentration threshold exceeded  $2.5 \text{ mg m}^{-3}$  and two higher thresholds were  
561 also considered. Sample sizes dwindle rapidly with increasing threshold. Of specific interest were the  
562 similarities and differences in respective bloom object sizes, their geographical location and collocation  
563 and timing. For the timing component the onset, duration and demise of individual bloom objects  
564 (events) could be considered. For the season all the identified space-time objects provided an estimate  
565 of the onset, duration and end of the bloom season as a whole. The season was found to be of similar  
566 length, but the onset was found to begin 26 days later in the AMM7v11 analyses than in the L4 product,  
567 and the AMM7v11 analyses persist the season for almost a month beyond the diagnosed end identified  
568 in the L4 product. Using traditional verification methods, data assimilation has been shown to  
569 considerably reduce the delay in bloom onset in the model (Skákala et al., 2020). Using feature-based  
570 verification methods, this study suggests that a substantial delay still remains.

571

572 There is a modest concentration bias in the AMM7v11 analyses compared to the L4 satellite ocean  
573 colour product. In this study we chose not to mitigate against this bias as it was not considered to  
574 impede the identification of bloom objects, which would prevent the ability of the methodology to  
575 identify matches and create paired object statistics. Any concentration bias does affect the results and  
576 this effect must be understood or at least kept in mind when interpreting results, in this case it will have  
577 contributed to the result that the AMM7v11 bloom objects are generally larger. An alternative approach  
578 would be to mitigate against the impact of the bias before using a threshold-based methodology such as  
579 MODE or MTD. A quantile mapping approach is available within the MODE tool (not yet available in  
580 MTD but should be available at some point) to remove the biases between two distributions as each  
581 temporal data set is analysed. Using this method the one threshold is fixed and the other threshold varies  
582 day-to-day (as shown in Figure 4). Another approach would be to analyse the bias for the whole season  
583 (as shown in Figure 3) and deriving an equivalent threshold from this larger data set, thus applying a

584 fixed threshold to all the days in the season, though there would still be two different thresholds applied  
585 to the two data sets.

586

587 MODE results suggest that the AMM7v11 bloom objects are larger than those in the L4 product.  
588 AMM7v11 produces more objects (in number) than seen in the L4 ocean colour product, yet many of  
589 the coastal objects seen in the L4 product are not as well resolved in AMM7v11 due to the coarseness of  
590 the coastline in the 7 km model. The additional AMM7v11 objects are mainly found in deeper Atlantic  
591 waters. The diagnosis of coastal blooms should improve if the model resolution were increased from  
592 7 km to 1.5 km.

593

594 Other work that formed part of this study, but is not reported on here, showed that constraining the Chl-  
595 *a* using assimilation of the satellite observations appears to benefit the model in terms of fewer  
596 unmatched bloom regions. This should translate to an improvement in the forecasts generated from this  
597 analysis compared with previous versions of the operational system and will be the subject of future  
598 work.

## 599 **6. Code availability**

600 Model Evaluation Tools (MET) was initially developed at the National Center for Atmospheric  
601 Research (NCAR) through grants from the National Science Foundation (NSF), the National Oceanic  
602 and Atmospheric Administration (NOAA), the United States Air Force (USAF) and the United States  
603 Department of Energy (DOE). The tool is now open source and available for download on github:  
604 <https://github.com/dtcenter/MET>. For this study MET version 8.1 of the software was used. MET  
605 allows for a variety of input file formats but some pre-processing of the CMEMS NetCDF files was  
606 necessary before the MODE package could be applied. This includes regridding of the observations  
607 onto the model grid, and addition of the forecast reference time variables to the NetCDF attributes. All  
608 aspects on the use of MET are provided in in the MET software documentation available online at  
609 <https://dtcenter.github.io/MET>.

## 610 **7. Data availability**

611 Data used in this paper was downloaded from the Copernicus Marine and Environment Monitoring  
612 Service (CMEMS). The datasets used were:

- 613 • [https://resources.marine.copernicus.eu/?option=com\\_csw&task=results?option=com\\_csw&view=de](https://resources.marine.copernicus.eu/?option=com_csw&task=results?option=com_csw&view=details&product_id=OCEANCOLOUR_ATL_CHL_L4_NRT_OBSERVATIONS_009_037)  
614 [tails&product\\_id=OCEANCOLOUR\\_ATL\\_CHL\\_L4\\_NRT\\_OBSERVATIONS\\_009\\_037](https://resources.marine.copernicus.eu/?option=com_csw&task=results?option=com_csw&view=details&product_id=OCEANCOLOUR_ATL_CHL_L4_NRT_OBSERVATIONS_009_037) (last  
615 access: August 2019),
- 616 • [https://resources.marine.copernicus.eu/?option=com\\_csw&view=details&product\\_id=NORTHWES](https://resources.marine.copernicus.eu/?option=com_csw&view=details&product_id=NORTHWESTSHELF_ANALYSIS_FORECAST_BIO_004_002_b)  
617 [TSHELF\\_ANALYSIS\\_FORECAST\\_BIO\\_004\\_002\\_b](https://resources.marine.copernicus.eu/?option=com_csw&view=details&product_id=NORTHWESTSHELF_ANALYSIS_FORECAST_BIO_004_002_b) (last access: August 2019)

618

619 The AMM7v11 analyses were not operational at the time of this study and not yet available from the  
620 CMEMS server.

## 621 **8. Author contribution**

622 All authors contributed to the introduction, data and methods, and conclusions. MM, RN, JM and CP  
623 contributed to the scientific evaluation and analysis of the results. MM and RN designed and ran the  
624 model assessments. CP supported the assessments through the provision and reformatting of the data  
625 used. DF provided detail on the model configurations used.

## 626 **9. Competing interests**

627 The authors declare that they have no conflict of interest.  
628

## 629 **10. Acknowledgements**

630 This study has been conducted using E.U. Copernicus Marine Service Information.

631

632 This work has been carried out as part of the Copernicus Marine Environment Monitoring Service  
633 (CMEMS) HiVE project. CMEMS is implemented by Mercator Ocean International in the framework  
634 of a delegation agreement with the European Union.

635

636 We would like to thank the National Center for Atmospheric Research (NCAR) Developmental Testbed  
637 Center (DTC) for the help received via their met\_help facility in getting MET to work with ocean data.

## 638 **11. References**

639 Allen, J. I. and Somerfield, P. J.: A multivariate approach to model skill assessment, *J. Mar. Syst.*,  
640 76(1–2), doi:10.1016/j.jmarsys.2008.05.009, 2009.

641 Allen, J. I., Holt, J. T., Blackford, J. and Proctor, R.: Error quantification of a high-resolution coupled  
642 hydrodynamic-ecosystem coastal-ocean model: Part 2. Chlorophyll-a, nutrients and SPM, *J. Mar. Syst.*,  
643 68(3–4), doi:10.1016/j.jmarsys.2007.01.005, 2007a.

644 Allen, J. I., Somerfield, P. J. and Gilbert, F. J.: Quantifying uncertainty in high-resolution coupled  
645 hydrodynamic-ecosystem models, *J. Mar. Syst.*, 64(1–4), doi:10.1016/j.jmarsys.2006.02.010, 2007b.

646 Antoine, D., Andrt, J. M. and Morel, A.: Oceanic primary production: 2. Estimation at global scale from  
647 satellite (Coastal Zone Color Scanner) chlorophyll, *Global Biogeochem. Cycles*, 10(1),  
648 doi:10.1029/95GB02832, 1996.

649 Anugerahanti, P., Roy, S. and Haines, K.: A perturbed biogeochemistry model ensemble evaluated  
650 against in situ and satellite observations, *Biogeosciences Discuss.*, doi:10.5194/bg-2018-136, 2018.

651 Behrenfeld, M. J., Boss, E., Siegel, D. A. and Shea, D. M.: Carbon-based ocean productivity and  
652 phytoplankton physiology from space, *Global Biogeochem. Cycles*, 19(1), doi:10.1029/2004GB002299,  
653 2005.

654 Bruggeman, J. and Bolding, K.: A general framework for aquatic biogeochemical models, *Environ.*  
655 *Model. Softw.*, 61, doi:10.1016/j.envsoft.2014.04.002, 2014.

656 Butenschön, M., Clark, J., Aldridge, J. N., Icarus Allen, J., Artioli, Y., Blackford, J., Bruggeman, J.,  
657 Cazenave, P., Ciavatta, S., Kay, S., Lessin, G., Van Leeuwen, S., Van Der Molen, J., De Mora, L.,  
658 Polimene, L., Saille, S., Stephens, N. and Torres, R.: ERSEM 15.06: A generic model for marine



659 biogeochemistry and the ecosystem dynamics of the lower trophic levels, *Geosci. Model Dev.*, 9(4),  
660 doi:10.5194/gmd-9-1293-2016, 2016.

661 Chelton, D. B., Schlax, M. G. and Samelson, R. M.: Global observations of nonlinear mesoscale eddies,  
662 *Prog. Oceanogr.*, 91(2), doi:10.1016/j.pocean.2011.01.002, 2011.

663 Chiswell, S. M.: Annual cycles and spring blooms in phytoplankton: Don't abandon Sverdrup  
664 completely, *Mar. Ecol. Prog. Ser.*, 443, doi:10.3354/meps09453, 2011.

665 Clark, A. J., Bullock, R. G., Jensen, T. L., Xue, M. and Kong, F.: Application of object-based time-  
666 domain diagnostics for tracking precipitation systems in convection-allowing models, *Weather*  
667 *Forecast.*, 29(3), doi:10.1175/WAF-D-13-00098.1, 2014.

668 Crocker, R., Maksymczuk, J., Mittermaier, M., Tonani, M. and Pequignet, C.: An approach to the  
669 verification of high-resolution ocean models using spatial methods, *Ocean Sci.*, 16(4), doi:10.5194/os-  
670 16-831-2020, 2020.

671 Crocker, R. L. and Mittermaier, M. P.: Exploratory use of a satellite cloud mask to verify {NWP}  
672 models, *Meteorol. Appl.*, 20, 197–205, 2013.

673 Davis, C., Brown, B. and Bullock, R.: Object-based verification of precipitation forecasts, Part {I}:  
674 Methods and application to mesoscale rain areas, *Mon. Wea. Rev.*, 134, 1772–1784, 2006.

675 Dorninger, M., Gilleland, E., Casati, B., Mittermaier, M., Ebert, E., Brown, B. and Wilson, L.: The set-  
676 up of the {M}esoscale {V}erification{I}nter-Comparison over {C}omplex {T}errain ({M}eso{VICT})  
677 project, *Bull. Amer. Meteorol. Soc.*, 2018.

678 Dutkiewicz, S., Hickman, A. E. and Jahn, O.: Modelling ocean-colour-derived chlorophyll A,  
679 *Biogeosciences*, 15(2), doi:10.5194/bg-15-613-2018, 2018.

680 Edwards, K. P., Barciela, R. and Butenschön, M.: Validation of the NEMO-ERSEM operational  
681 ecosystem model for the North West European continental shelf, *Ocean Sci.*, 8(6), doi:10.5194/os-8-  
682 983-2012, 2012.

683 Falkowski, P. G., Barber, R. T. and Smetacek, V.: Biogeochemical controls and feedbacks on ocean  
684 primary production, *Science* (80-. ), 281(5374), doi:10.1126/science.281.5374.200, 1998.

685 Ford, D. A., Van Der Molen, J., Hyder, K., Bacon, J., Barciela, R., Creach, V., McEwan, R., Ruardij, P.  
686 and Forster, R.: Observing and modelling phytoplankton community structure in the North Sea,

687 Biogeosciences, 14(6), doi:10.5194/bg-14-1419-2017, 2017.

688 Gilleland, E., Ahijevych, D., Brown, B. and Ebert, E.: Intercomparison of Spatial Forecast Verification  
689 Methods, *Wea. Forecast.*, 24, 2009.

690 Gilleland, E., Lindström, J. and Lindgren, F.: Analyzing the image warp forecast verification method on  
691 precipitation fields from the {ICP}, *Weather Forecast.*, 25(4), 1249–1262, 2010.

692 Gordon, H. R., Clark, D. K., Brown, J. W., Brown, O. B., Evans, R. H. and Broenkow, W. W.:  
693 Phytoplankton pigment concentrations in the Middle Atlantic Bight: comparison of ship determinations  
694 and CZCS estimates, *Appl. Opt.*, 22(1), doi:10.1364/ao.22.000020, 1983.

695 Hausmann, U. and Czaja, A.: The observed signature of mesoscale eddies in sea surface temperature  
696 and the associated heat transport, *Deep. Res. Part I Oceanogr. Res. Pap.*, 70,  
697 doi:10.1016/j.dsr.2012.08.005, 2012.

698 Hipsey, M. R., Gal, G., Arhonditsis, G. B., Carey, C. C., Elliott, J. A., Frassl, M. A., Janse, J. H., de  
699 Mora, L. and Robson, B. J.: A system of metrics for the assessment and improvement of aquatic  
700 ecosystem models, *Environ. Model. Softw.*, 128, doi:10.1016/j.envsoft.2020.104697, 2020.

701 Jolliff, J. K., Kindle, J. C., Shulman, I., Penta, B., Friedrichs, M. A. M., Helber, R. and Arnone, R. A.:  
702 Summary diagrams for coupled hydrodynamic-ecosystem model skill assessment, *J. Mar. Sys.*, 76, 64–  
703 82, 2009.

704 King, R. R., While, J., Martin, M. J., Lea, D. J., Lemieux-Dudon, B., Waters, J. and O’Dea, E.:  
705 Improving the initialisation of the Met Office operational shelf-seas model, *Ocean Model.*, 130,  
706 doi:10.1016/j.ocemod.2018.07.004, 2018.

707 LORENZEN, C. J.: SURFACE CHLOROPHYLL AS AN INDEX OF THE DEPTH, CHLOROPHYLL  
708 CONTENT, AND PRIMARY PRODUCTIVITY OF THE EUPHOTIC LAYER, *Limnol. Oceanogr.*,  
709 15(3), doi:10.4319/lo.1970.15.3.0479, 1970.

710 Madec, G. and the N. team: *Nemo Engine.*, 2016.

711 Mass, C. F., Ovens, D., Westrick, K. and Colle, B. A.: Does increasing horizontal resolution produce  
712 more skillful forecasts? The results of two years of real-time numerical weather prediction over the  
713 Pacific northwest, *Bull. Amer. Meteorol. Soc.*, 83(3), 407–430, 2002.

714 Mittermaier, M. and Bullock, R.: Using {MODE} to explore the spatial and temporal characteristics of

715 cloud cover forecasts from high-resolution {NWP} models, *Meteorol. Appl.*, 20, 187–196, 2013.

716 Mittermaier, M., North, R., Semple, A. and Bullock, R.: Feature-based diagnostic evaluation of global  
717 NWP forecasts, *Mon. Wea. Rev.*, 144(10), Submitted, 2016.

718 Moore, T. S., Campbell, J. W. and Dowell, M. D.: A class-based approach to characterizing and  
719 mapping the uncertainty of the MODIS ocean chlorophyll product, *Remote Sens. Environ.*, 113(11),  
720 2424–2430, doi:<https://doi.org/10.1016/j.rse.2009.07.016>, 2009.

721 De Mora, L., Butenschön, M. and Allen, J. I.: The assessment of a global marine ecosystem model on  
722 the basis of emergent properties and ecosystem function: A case study with ERSEM, *Geosci. Model  
723 Dev.*, 9(1), doi:10.5194/gmd-9-59-2016, 2016.

724 Morrow, R. and Le Traon, P. Y.: Recent advances in observing mesoscale ocean dynamics with satellite  
725 altimetry, *Adv. Sp. Res.*, 50(8), doi:10.1016/j.asr.2011.09.033, 2012.

726 O’Dea, E. J., Arnold, A. K., Edwards, K. P., Furner, R., Hyder, P., Martin, M. J., Siddorn, J. R.,  
727 Storkey, D., While, J., Holt, J. T. and Liu, H.: An operational ocean forecast system incorporating  
728 NEMO and SST data assimilation for the tidally driven European North-West shelf, *J. Oper. Oceanogr.*,  
729 5(1), doi:10.1080/1755876X.2012.11020128, 2012.

730 O’Dea, E., Furner, R., Wakelin, S., Siddorn, J., While, J., Sykes, P., King, R., Holt, J. and  
731 Hewitt, H.: The CO5 configuration of the 7&thinsp;km Atlantic Margin Model: Large scale biases  
732 and sensitivity to forcing, physics options and vertical resolution, *Geosci. Model Dev. Discuss.*,  
733 doi:10.5194/gmd-2017-15, 2017.

734 Rossa, A. M., Nurmi, P. and Ebert, E. E.: *Precipitation: Advances in Measurement, Estimation and  
735 Prediction*, pp. 418–450, Springer., 2008.

736 Saux Picart, S., Butenschén, M. and Shutler, J. D.: Wavelet-based spatial comparison technique for  
737 analysing and evaluating two-dimensional geophysical model fields, *Geosci. Model Dev.*, 5(1),  
738 doi:10.5194/gmd-5-223-2012, 2012.

739 Schalles, J. F.: Optical remote sensing techniques to estimate phytoplankton chlorophyll a  
740 concentrations in coastal waters with varying suspended matter and cdom concentrations, in *Remote  
741 Sensing and Digital Image Processing*, vol. 9., 2006.

742 Shutler, J. D., Smyth, T. J., Saux-Picart, S., Wakelin, S. L., Hyder, P., Orekhov, P., Grant, M. G.,

743 Tilstone, G. H. and Allen, J. I.: Evaluating the ability of a hydrodynamic ecosystem model to capture  
744 inter- and intra-annual spatial characteristics of chlorophyll-a in the north east Atlantic, *J. Mar. Syst.*,  
745 88(2), doi:10.1016/j.jmarsys.2011.03.013, 2011.

746 Siegel, D. A., Doney, S. C. and Yoder, J. A.: The North Atlantic Spring Phytoplankton Bloom and  
747 Sverdrup's Critical Depth Hypothesis, *Science* (80-. ), 296(5568), 730–733,  
748 doi:10.1126/science.1069174, 2002.

749 Skákala, J., Ford, D., Brewin, R. J. W., McEwan, R., Kay, S., Taylor, B., de Mora, L. and Ciavatta, S.:  
750 The Assimilation of Phytoplankton Functional Types for Operational Forecasting in the Northwest  
751 European Shelf, *J. Geophys. Res. Ocean.*, 123(8), 5230–5247, doi:10.1029/2018JC014153, 2018.

752 Skákala, J., Bruggeman, J., Brewin, R. J. W., Ford, D. A. and Ciavatta, S.: Improved Representation of  
753 Underwater Light Field and Its Impact on Ecosystem Dynamics: A Study in the North Sea, *J. Geophys.*  
754 *Res. Ocean.*, 125(7), e2020JC016122, doi:10.1029/2020JC016122, 2020.

755 Smyth, T. J., Allen, I., Atkinson, A., Bruun, J. T., Harmer, R. A., Pingree, R. D., Widdicombe, C. E.  
756 and Somerfield, P. J.: Ocean net heat flux influences seasonal to interannual patterns of plankton  
757 abundance, *PLoS One*, 9(6), e98709, doi:10.1371/journal.pone.0098709, 2014.

758 Stow, C. A., Jolliff, J., McGillicuddy, D. J., Doney, S. C., Allen, J. I., Friedrichs, M. A. M., Rose, K. A.  
759 and Wallhead, P.: Skill assessment for coupled biological/physical models of marine systems, *J. Mar.*  
760 *Syst.*, 76(1–2), doi:10.1016/j.jmarsys.2008.03.011, 2009.

761 Sverdrup, H. U.: On conditions for the vernal blooming of phytoplankton, *ICES J. Mar. Sci.*, 18(3),  
762 doi:10.1093/icesjms/18.3.287, 1953.

763 Taylor, K. E.: Summarizing multiple aspects of model performance in a single diagram, *J. Geophys.*  
764 *Res. Atmos.*, 106(D7), doi:10.1029/2000JD900719, 2001.

765 Le Traon, P. Y., Reppucci, A., Fanjul, E. A., Aouf, L., Behrens, A., Belmonte, M., Bentamy, A.,  
766 Bertino, L., Brando, V. E., Kreiner, M. B., Benkiran, M., Carval, T., Ciliberti, S. A., Claustre, H.,  
767 Clementi, E., Coppini, G., Cossarini, G., De Alfonso Alonso-Muñoyerro, M., Delamarche, A.,  
768 Dibarboure, G., Dinessen, F., Drevillon, M., Drillet, Y., Faugere, Y., Fernández, V., Fleming, A.,  
769 Garcia-Hermosa, M. I., Sotillo, M. G., Garric, G., Gasparin, F., Giordan, C., Gehlen, M., Gregoire, M.  
770 L., Guinehut, S., Hamon, M., Harris, C., Hernandez, F., Hinkler, J. B., Hoyer, J., Karvonen, J., Kay, S.,

771 King, R., Lavergne, T., Lemieux-Dudon, B., Lima, L., Mao, C., Martin, M. J., Masina, S., Melet, A.,  
772 Nardelli, B. B., Nolan, G., Pascual, A., Pistoia, J., Palazov, A., Piolle, J. F., Pujol, M. I., Pequignet, A.  
773 C., Peneva, E., Gómez, B. P., de la Villeon, L. P., Pinardi, N., Pisano, A., Pouliquen, S., Reid, R.,  
774 Remy, E., Santoleri, R., Siddorn, J., She, J., Staneva, J., Stoffelen, A., Tonani, M., Vandenbulcke, L.,  
775 von Schuckmann, K., Volpe, G., Wettre, C. and Zacharioudaki, A.: From observation to information  
776 and users: The Copernicus Marine Service Perspective, *Front. Mar. Sci.*, 6(May),  
777 doi:10.3389/fmars.2019.234, 2019.

778 Vichi, M., Allen, J. I., Masina, S. and Hardman-Mountford, N. J.: The emergence of ocean  
779 biogeochemical provinces: A quantitative assessment and a diagnostic for model evaluation, *Global*  
780 *Biogeochem. Cycles*, 25(2), doi:10.1029/2010GB003867, 2011.

781 Waters, J., Lea, D. J., Martin, M. J., Mirouze, I., Weaver, A. and While, J.: Implementing a variational  
782 data assimilation system in an operational 1/4 degree global ocean model, *Q. J. R. Meteorol. Soc.*,  
783 141(687), 333–349, doi:10.1002/qj.2388, 2015.

784 Winder, M. and Cloern, J. E.: The annual cycles of phytoplankton biomass, *Philos. Trans. R. Soc. B*  
785 *Biol. Sci.*, 365(1555), doi:10.1098/rstb.2010.0125, 2010.

786

787