



# 1 Testing the reconstruction of modelled particulate organic 2 carbon from surface ecosystem components using PlankTOM12 3 and Machine Learning

4  
5 Anna Denvil-Sommer<sup>1</sup>, Erik T. Buitenhuis<sup>1</sup>, Rainer Kiko<sup>2,3</sup>, Fabien Lombard<sup>2,4</sup>, Lionel Guidi<sup>2</sup>, Corinne Le  
6 Quéré<sup>1</sup>

7  
8 <sup>1</sup>School of Environmental Science, University of East Anglia, Norwich, UK

9 <sup>2</sup>Sorbonne Université, Centre National de la Recherche Scientifique (CNRS), Laboratoire d'Océanographie de  
10 Villefranche (LOV), Villefranche-sur-Mer, France

11 <sup>3</sup>GEOMAR Helmholtz Center for Ocean Research, Kiel, Germany

12 <sup>4</sup>Institut Universitaire de France (IUF), Paris, France

13 *Correspondence to:* Anna Denvil-Sommer (anna.sommer.lab@gmail.com)

14 **Abstract.** Understanding the relationship between surface marine ecosystems and the export of carbon to depth by  
15 sinking organic particles is key to represent the effect of ecosystem dynamics and diversity, and their evolution under  
16 multiple stressors, on the carbon cycle and climate in models. Recent observational technologies have greatly  
17 increased the amount of data available, both for the abundance of diverse plankton groups and for the concentration  
18 and properties of particulate organic carbon in the ocean interior. Here we use synthetic model data to test the potential  
19 of using Machine Learning (ML) to reproduce concentrations of particulate organic carbon within the ocean interior  
20 based on surface ecosystem and environmental data. We test two machine learning methods that differ in their  
21 approaches to data-fitting, the Random Forest and XGBoost methods. The synthetic data is sampled from the  
22 PlankTOM12 global biogeochemical model using the time and coordinates of existing observations. We test 27  
23 different combinations of possible drivers to reconstruct small (POC<sub>S</sub>) and large (POC<sub>L</sub>) particulate organic carbon  
24 concentrations. We show that ML can successfully be used to reproduce modelled particulate organic carbon over  
25 most of the ocean based on ecosystem and modelled environmental drivers. XGBoost showed better results compared  
26 to Random Forest thanks to its gradient boosting trees architecture. The inclusion of Plankton Functional Types (PFTs)  
27 in driver sets improved the accuracy of the model reconstruction by 58% on average for POC<sub>S</sub>, and by 22% for POC<sub>L</sub>.  
28 Results were less robust over the Equatorial Pacific and some parts of the high latitudes. For POC<sub>S</sub> reconstruction, the  
29 most important drivers were the depth level, temperature, microzooplankton and PO<sub>4</sub>, while for POC<sub>L</sub> it was the depth  
30 level, temperature, mixed-layer depth, microzooplankton, phaeocystis, PO<sub>4</sub> and chlorophyll *a* averaged over the  
31 mixed-layer depth. These results suggest that it will be possible to identify linkages between surface environmental  
32 and ecosystem structure and particulate organic carbon distribution within the ocean interior using real observations,  
33 and to use this knowledge to improve both our understanding of ecosystem dynamics and of their functional  
34 representation within models.

## 35 36 1. Introduction.

37  
38 Progress in numerical ocean modelling over multiple decades coupled with fundamental knowledge of fluid dynamics  
39 have led to an explicit representation of ocean dynamics in Earth System Models and of most of its key features, apart  
40 from small-scale features which are parametrized. In contrast, ecosystem dynamics in ocean biogeochemical models  
41 are much more reliant on empirical data for growth and loss processes, with the theoretical basis limited to the dynamic  
42 representation of interactions among lower trophic levels (zooplankton and smaller organisms) and their influence on  
43 carbon pools and fluxes (Le Quéré et al., 2005; Hood et al., 2006). The recent advances in observational technologies  
44 including imaging data (Guidi et al., 2016), genomics (Kirchman et al., 2016), and field study (Mutshinda et al., 2017;  
45 Batten et al., 2019, Lombard et al., 2019), offer new opportunities to improve our understanding of marine ecosystem  
46 dynamics, and to better represent its influence on carbon pools and fluxes in models that are used to project future  
47 climate change and associated impacts on ecosystems.

48 One strategy to represent lower trophic interactions in global biogeochemical models is to combine different species  
49 into Plankton Functional Types (PFTs) based on their unique influence on global biogeochemical cycles (Le Quéré et  
50 al., 2005; Hood et al., 2006). This approach enables the representation of plankton types that are unique, have an



51 influence on other PFTs within the ecosystem and are of quantitative importance for carbon flux and other  
52 biogeochemical fluxes. The PlankTOM12 model is among the most detailed in this category of models with its  
53 inclusion of an explicit representation of twelve PFTs: six phytoplankton, five zooplankton, and bacteria.  
54 PlankTOM12 builds on the published version PlankTOM10 (Le Quéré et al., 2016) that has been extended to include  
55 gelatinous zooplankton (Wright et al., 2021) and pteropods (Buitenhuis et al., 2019). Much effort has been put into  
56 the development of PFTs and associated representation of surface ecosystem dynamics, which has led to the  
57 demonstration that: (1) the representation of trophic levels was a key determinant of the low chlorophyll *a*  
58 concentration observed in the Southern Ocean summer (Le Quéré et al., 2016); (2) CaCO<sub>3</sub> dissolution above the  
59 lysocline is needed to reproduce observations of both biomass and export of PFT calcifiers, and (3) gelatinous  
60 zooplankton plays an important role in determining surface biomass of other PFTs (Wright et al., 2021).

61 In contrast, the transfer of organic matter resulting from surface ecosystem dynamics into carbon exported to the deep  
62 ocean via the sinking of particulate organic matter has received much less attention, so that improvements in the  
63 representation of the PFTs do not necessarily translate into improvements in sinking of particulate matter (Wright et  
64 al., 2021). The export flux of particulate organic carbon from the surface ocean to depth is around 10 PgC yr<sup>-1</sup>  
65 (Schlitzer, 2002), which is as large as the CO<sub>2</sub> emitted to the atmosphere by human activities and nearly four times  
66 larger than the mean oceanic CO<sub>2</sub> sink in recent decades (Friedlingstein et al., 2022). Changes in carbon exported to  
67 depth can have a large impact on air-sea CO<sub>2</sub> fluxes and on the amount of CO<sub>2</sub> emissions that remain in the atmosphere  
68 where they cause climate change.

69 The growing amount of observations provides the opportunity to develop a new approach to explore the linkages  
70 between surface ecosystem dynamics and the distribution of particulate organic carbon in the ocean, and to improve  
71 the representation of particle sinking fluxes in models. However, there is a risk of over-interpreting the data by  
72 applying Machine Learning (ML) methods directly to link the observed surface environment and ecosystem structure  
73 with the observed particulate organic carbon distribution. The use of synthetic observations based on model data  
74 therefore provides a minimum test to assess the likely success and usefulness of such an approach.

75 ML has been widely used in biogeochemical and geophysical applications and provided efficient results in  
76 reconstructions of ocean surface pCO<sub>2</sub> (Friedrich and Oschlies, 2009; Telszewski et al., 2009; Landschützer et al.,  
77 2013; Denvil-Sommer et al., 2019) and of particulate organic carbon (Sauzède et al., 2016, 2017) as well as in the  
78 analysis of driver importance (Sauzède et al., 2020).

79 Here we use model data to verify the hypothesis that the composition of surface ecosystems and environmental  
80 conditions are indeed reflected in the abundance and size of the organic particles in the ocean interior. We reconstruct  
81 the concentration of organic particles as represented by small (POC<sub>s</sub>, particles < 256µm) and large (POC<sub>L</sub>, particles >  
82 256µm) particulate carbon in the PlankTOM12 model. Using this information alongside with modelled environmental  
83 and ecosystem conditions we develop a ML method to reproduce POC<sub>s</sub> and POC<sub>L</sub> over the global ocean and verify  
84 the hypothesis. This constitutes a necessary although not sufficient test that the approach can subsequently be used to  
85 reveal linkages using real observations and to inform model developments.

## 86 2. Data and Methods.

87 In this section we describe a set of variables that will be used to test the ML method's ability to reconstruct particulate  
88 organic carbon concentrations based on ocean model data. We create a set of synthetic data by sampling a model at  
89 the time and location of real-world observations. We discuss the availability and distribution of real-world  
90 observations and their limitations. In this section we also describe the PlankTOM12 global ocean biogeochemical  
91 model and how we use it to develop a ML method and test its ability to reconstruct small and large particulate organic  
92 carbon with a limited number of observations. To provide resemblance to the real data availability we focus on the  
93 period 2009-2013 which guarantees additional sampling of co-located biological, chemical, and environmental  
94 variables from the Tara expeditions (Sunagawa et al., 2020).

95 Two sets of data are needed to test the Machine Learning method: a set of targets and a set of drivers. The drivers  
96 represent the input variables to the ML method (here the biological, chemical, and environmental variables). The  
97 targets represent the variables we are trying to reconstruct (here the particulate organic matter POC<sub>s</sub> and POC<sub>L</sub>). The  
98 ML will then determine the relationship between the drivers and targets, which can then be applied in regions where  
99 drivers are available to infer targets where the later data do not exist.

100

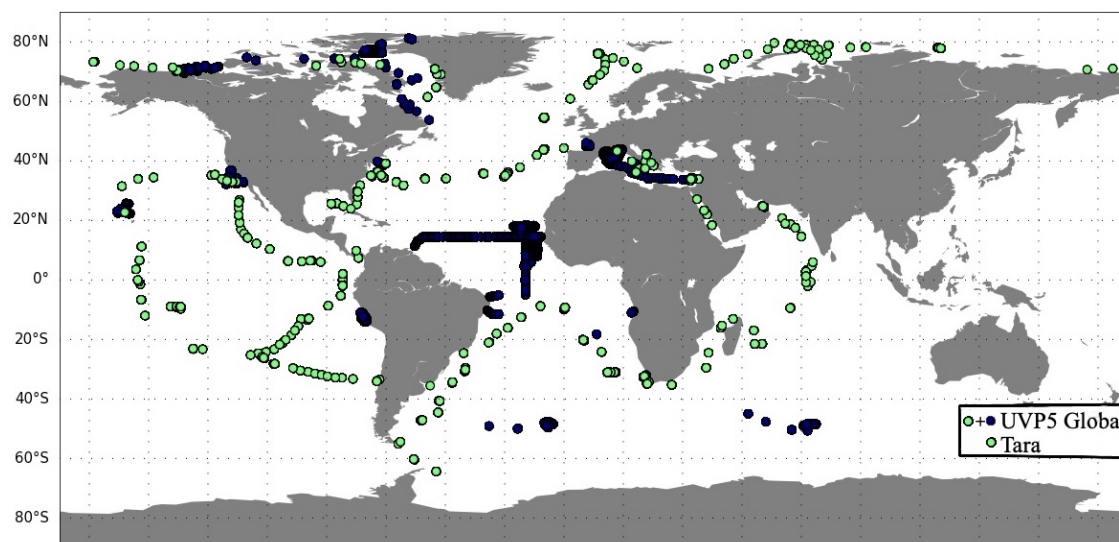


101 **2.1.1. Measurements of particle size distributions and concentrations (the targets).**

102 We use observations of particle distribution in two ways. First to determine the time and location of the observations,  
103 and second to verify that the ocean model is of sufficient quality to be used in this analysis. The sampling of the  
104 particulate organic carbon concentration is based on the data from an Underwater Vision Profiler 5 (UVP5) (Gorsky  
105 et al., 2000, 1992; Picheral et al., 2010; Kiko et al., 2022). UVP5 measures particles of size from 50  $\mu\text{m}$  to a few mm.  
106 For the purpose of comparing the UVP5 data with the PlankTOM12 model data, we converted measured biovolume  
107 concentration ( $\text{mm}^3/\text{L}$ ) of particles to carbon biomass concentrations ( $\mu\text{mol}/\text{L}$ ) using the empirical equation from  
108 Alldredge (1998) for particulate organic carbon:

$$\text{BM} (\mu\text{g}) = 0.99 * \text{BV} (\text{mm}^3)^{0.52}$$

109  
110



111  
112 **Figure 1. Location of the observations from the UVP5 database over the period 2009-2013. Green dots correspond to Tara**  
113 **expeditions, and were included in the global UVP5 database.**

114 We summed size classes from 50.8  $\mu\text{m}$  to 256  $\mu\text{m}$  for the small particulate organic carbon (POC<sub>S</sub>) and from 256  $\mu\text{m}$   
115 to 5.16 mm for large particulate organic carbon (POC<sub>L</sub>). POCs below 100  $\mu\text{m}$  are not well captured by the UVP sensor,  
116 which therefore underestimates this size-class of aggregated particles. We extrapolated the total size of particles up to  
117 0.001mm by using the size spectra theory to provide a better estimate of POC biomass concentration in line with the  
118 model. Following Guidi et al. (2008) we used the abundance of particles sized from 0.250 to 1.5 mm excluding rare  
119 particles to estimate the coefficients of logarithmic relationship between the size of particles and its abundance:

$$\log(\text{abundance}) = a * \log(\text{size}) + b$$

120  
121 Using this equation we estimated the abundance of particles of size less than 100  $\mu\text{m}$ .

122 There are 2603 vertical profiles of UVP5 measurements during 2009-2013, including 752 profiles which are co-  
123 located with the stations from the Tara expeditions that provide the environmental and ecosystem variables (Figure 1;  
124 Section 2.1.2). The measurements are sparse in time and space. There are no measurements in the Southern Ocean,  
125 Western Pacific Ocean and Eastern Indian Oceans.

126 **2.1.2. Measurements of environmental and ecosystem variables (the drivers).**

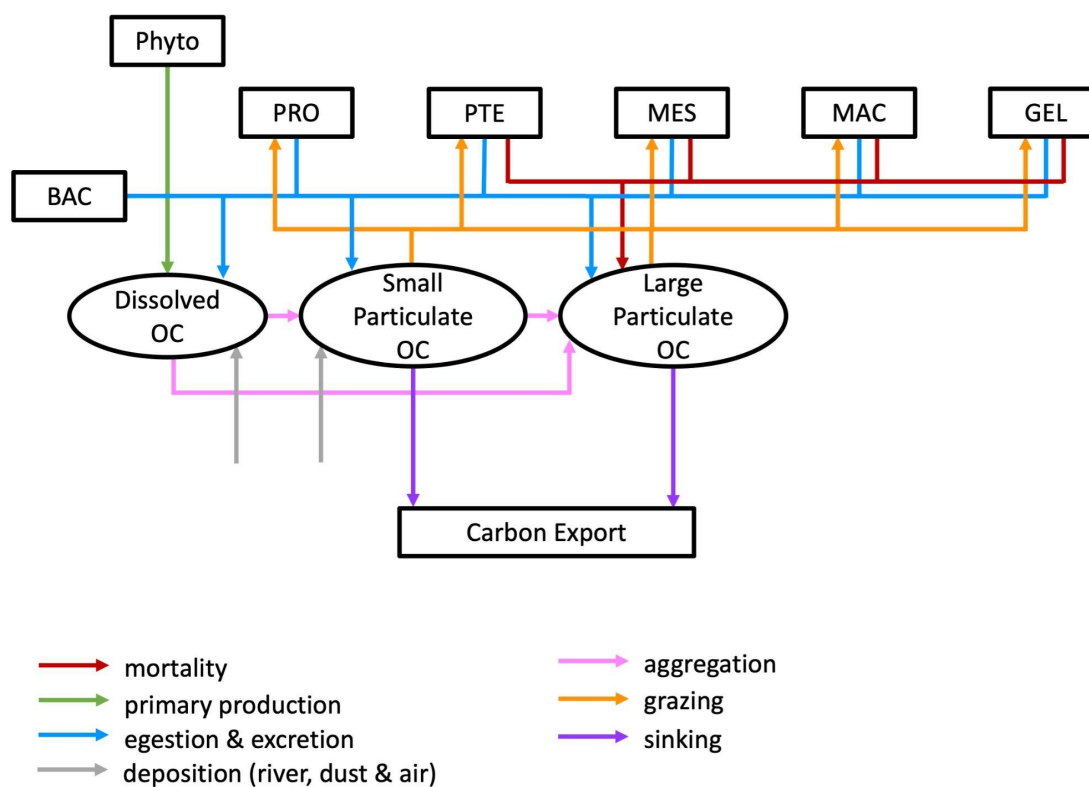
127 We use observations of environmental and ecosystem variables to determine the time and location of the observations  
128 that are collocated with the target variables. To represent the main physical and chemical drivers responsible for the  
129 concentration and variability of POC<sub>S</sub> and POC<sub>L</sub> we use measurements of ocean temperature, chlorophyll *a*, phosphate  
130 PO<sub>4</sub>, nitrates NO<sub>3</sub>, mixed-layer depth (MLD). These variables were measured during Tara expeditions along with the



131 particle size distributions and concentrations using UVP instruments onboard these cruises. However, chlorophyll *a*,  
 132 PO<sub>4</sub>, and nitrates were not measured systematically at each depth level. Thus, their averages over MLD are tested as  
 133 possible drivers as well. To represent the biological drivers, we use information on PFTs.

134 **2.1.3. The NEMO-PlankTOM12 Global biogeochemical model.**

135 We used the output from the NEMO-PlankTOM12 coupled physical–biogeochemical model of the global ocean at  
 136 daily and monthly time resolution. NEMO represents physical transport processes and is used in its v3.6-ORCA2  
 137 version, with a horizontal resolution 2° longitude and 0.3° to 1.5° latitude, and 31 vertical levels. It is forced by daily  
 138 meteorological data from NCEP reanalysis (Kalnay et al., 1996) over the period 1948–2020, with output for 2009–  
 139 2013 used here. This model version is identical to that used to estimate the ocean CO<sub>2</sub> sink in the Global Carbon  
 140 Budget 2021 annual update (Friedlingstein et al. 2021).







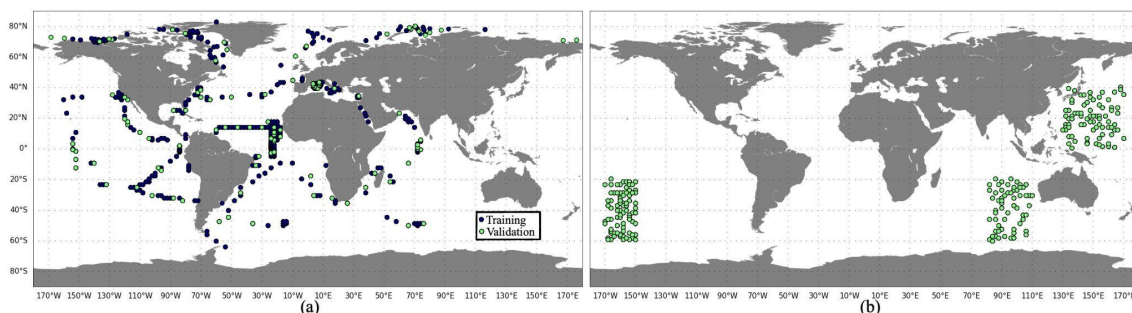
152 PlankTOM12 represents sinking processes through the explicit representation of two organic particle of different size,  
153 with small particles sinking at a constant speed of 3 m/d, and larger particles sinking at a variable speed between 3  
154 and 150 m/d depending on the ballast effect of their mineral content (Buitenhuis et al., 2013). In addition, a dissolved  
155 organic carbon component is transported via ocean currents. Particles are generated through mass flux from the PFTs  
156 resulting from mortality and egestion and from aggregation through differential sinking or turbulent coagulation, and  
157 destroyed through grazing by zooplankton and remineralisation by bacteria and through disaggregation from shear  
158 currents. Large PFTs contribute mostly to  $POC_L$ , while small PFTs contribute mostly to  $POC_S$ . (Le Quéré et al. 2016;  
159 Fig. 2).

160 The NEMO-PlankTOM12 model output was sampled at the time and location identified from the observations  
161 mentioned above to create a synthetic dataset. The model grid-coordinate closest to the real geographical position was  
162 chosen. If several measurements were co-localised at the same grid coordinate and same time step (day for daily  
163 PlankTOM12 and month for monthly PlankTOM12 outputs), it is counted as one measurement. This model sampling  
164 produced 400 positions when using the daily or monthly PlankTOM12 outputs. All drivers and targets were taken  
165 from the model output at the corresponding coordinates up to 1400 m depth. These outputs served as the reference for  
166 validation and evaluation of the ML methods and for establishing the sets of the most important drivers.

## 167 **2.2. Method.**

168 We tested 2 ML methods that are widely used in target's reconstruction based on tabular data sets: the Random Forest  
169 regressor and the XGBoost (Extreme Gradient Boosting) regressor. The Random Forest (RF) regressor is an ensemble  
170 algorithm that contains a number of decision trees on various subsets of the given dataset and takes as output the  
171 average of prediction from each tree estimator. RF can run several trees at the same time allowing a use of a large  
172 number of input variables, and it is robust to overfitting (Biau, 2012). XGBoost (XGB) regressor is an effective tree-  
173 based ensemble learning algorithm (Chen and Guestrin, 2016). It builds several models sequentially where each new  
174 model attempts to correct errors from the previous one. XGBoost uses the gradient descent algorithm to minimise the  
175 loss function of the model. Using RF and XGBoost we can estimate the driver importance to identify which driver has  
176 the greatest impact on the predictions. To check the driver importance, we use `drop_col_feat_imp` python function  
177 (<https://gist.github.com/erykml/6854134220276b1a50862aa486a44192>). This method estimates how the accuracy of  
178 the ML output changes if one of the drivers is dropped off from a driver set (DS) based on the training dataset.

179 Effective ML algorithm requires sets of training, validation and test data. The training data builds up the ML model.  
180 Model evaluates training data repeatedly to learn about the relationship between inputs (driver set) and known outputs  
181 (target set) and adjusts itself to better represent the target. The purpose of validation data is to evaluate the model  
182 during its training by introducing new unseen data. It allows us to evaluate how a developed model works on a new  
183 dataset and to optimise hyperparameters. The test data evaluate the final accuracy of the ML model and confirm that  
184 the model works correctly on any unseen data. It is new data that did not participate in the training algorithm. The  
185 accuracy is worse on validation and test data compared to training data set. The difference in model performance on  
186 training and validation data can signal an overfitting, while this difference between validation and test data can  
187 demonstrate an effect of data mismatch. It is worth noting that RF does not necessarily need validation data set as they  
188 perform internal validation. During the training algorithm each tree is constructed from a random subset of original  
189 data, usually it represents two thirds of data and one third of data is used to estimate out-of bag error to assess model  
190 performance. XGB uses a validation data set to evaluate the model during training and to prevent overfitting by  
191 applying an early stopping. In the present study the available data were split into training and validation data sets (Fig.  
192 3a). Validation data is not included in RF training, however we use it to test the performance of trained RF and tune  
193 hyperparameters afterwards. The test data are taken from the regions where there are no observations (Fig. 3b): 3  
194 months for each year from the period 2009-2013 and 6 positions for each month were chosen randomly. This will  
195 allow us to identify the possible accuracy of reconstruction that can be reached in these regions when we will apply a  
196 developed method to real observations. However, when  $POC_S$  and  $POC_L$  will be reconstructed using only real-world  
197 observations, we will need to split all available data into training, validation and test data sets.



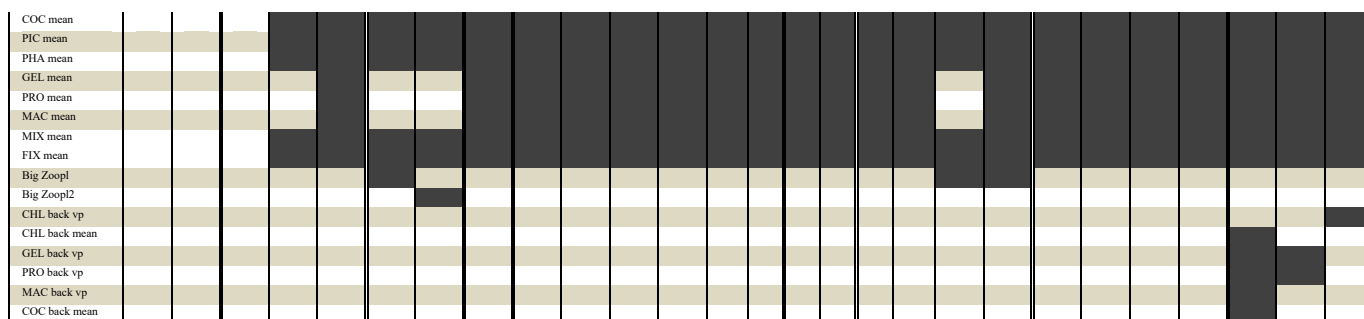
198  
 199 **Figure 3. The spatial distribution of: (a) - training (blue) and validation (green) data sets; (b) - test data set; based on**  
 200 **PlankTOM12 monthly outputs.**

201 We use RandomForestRegressor function from scikit-learn ([https://scikit-](https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html)  
 202 [learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html](https://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html)) with its default parameters and  
 203 min\_sample\_leaf equals 20. To apply XGBoost regressor we use XGBRegressor from xgboost  
 204 ([https://xgboost.readthedocs.io/en/stable/python/python\\_intro.html](https://xgboost.readthedocs.io/en/stable/python/python_intro.html)). Parameters were set as follows  
 205 n\_estimators=2000, max\_depth=7, eta=0.01, subsample=0.7, colsample\_bytree=0.8, gamma=0.01 for POC<sub>L</sub> and  
 206 gamma= 0.3 for POC<sub>S</sub>, early\_stopping\_rounds = 10.

207 We tested 27 driver sets (DSs) that are summarised in Table 1. For each DS we identify the most important drivers  
 208 that influenced the reconstruction of small (POC<sub>S</sub>) and large (POC<sub>L</sub>) particulate organic carbon concentration. The  
 209 drivers include geographic variables (depth, sin(latitude), cos(longitude)), physical variables (incident light, MLD,  
 210 co-located temperature), chemical variables (PO<sub>4</sub>, NO<sub>3</sub>, including co-located values and averages over the MLD), and  
 211 biological variables (chlorophyll *a*, 12 PFTs listed above: DIA, MIX, COC, PIC, PHA, FIX, PRO, PTE, MES, GEL,  
 212 BAC, including co-located values and averages over the MLD).

213 **Table 1. Compounds of driver's sets: dark grey cells correspond to the drivers present in the driver set. 'vp' – vertical**  
 214 **profile, 'mean' – average over MLD, 'back' – values from previous month.**

Driver set	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20	21	22	23	24	25	26	27
Drivers																											
depth																											
Sin(lat)																											
Sin(long)																											
Cos(long)																											
Incident light																											
MLD																											
Temperature vp																											
CHL vp																											
CHL mean																											
NO <sub>3</sub> vp																											
PO <sub>4</sub> vp																											
NO <sub>3</sub> mean																											
PO <sub>4</sub> mean																											
BAC vp																											
MES vp																											
PTE vp																											
DIA vp																											
COC vp																											
PIC vp																											
PHA vp																											
GEL vp																											
PRO vp																											
MAC vp																											
MIX vp																											
FIX vp																											
BAC mean																											
MES mean																											
PTE mean																											
DIA mean																											



215

216 The driver sets can be split into 9 thematic groups which together test the role of PFTs and sub-classes within, the role  
 217 of surface versus depth profiles for some variables, and the role of information from the previous month:

- 218 I. No PFTs (short name (sh.n.) ‘No PFT’): Driver sets 1 and 2 do not include any PFTs and focus on the  
 219 influence of temperature, MLD, chlorophyll *a*, NO<sub>3</sub> and PO<sub>4</sub> on POC<sub>s</sub> and POC<sub>L</sub> reconstruction.
- 220 II. Introduction of PFTs (sh.n. ‘PFT introduction’): DSs 3, 4 and 5 are dedicated to the investigation of the  
 221 introduction of PFTs in the reconstruction. In DS 3 we introduced 12 PFTs vertical profiles, even though this  
 222 information will be challenging to reproduce with observations due to the lack of data. Nevertheless, it is  
 223 important to test the capacity of ML if all 12 PFTs were available over the depth. DS 4 includes the vertical  
 224 profiles of 6 heterotrophs (zooplanktons and bacteria) because they contribute to influencing the vertical  
 225 distribution of POC<sub>s</sub> and POC<sub>L</sub>, and 6 phytoplankton averaged over MLD because they are responsible for  
 226 primary production. In DS 5 we added averages over MLD of the 6 heterotrophs that were not included in  
 227 DS 4.
- 228 III. Big zooplankton (sh.n. ‘Zooplankton combined’): In DSs 6 and 7 we tested the influence of big zooplanktons  
 229 summed into one variable to account for their combined effect rather than the distinctions among PFTs. The  
 230 big zooplankton is represented by the sum of mesozooplankton, gelatinous zooplankton and  
 231 macrozooplankton in DS 6, with the addition of pteropod in DS 7.
- 232 IV. Exclusion of bacteria (sh.n. ‘No vertical BAC’): DS 8 does not have a bacteria (BAC) vertical profile  
 233 compared to set 5.
- 234 V. Individual zooplankton types (sh.n. ‘Individual PFT’): DSs 9, 10, 11, 12, 13 and 14 test the influence of  
 235 individual types of heterotrophs, bacteria (BAC), microzooplankton (PRO), pteropod (PTE),  
 236 mesozooplankton (MES), gelatinous zooplankton (GEL), microzooplankton (MAC), respectively.
- 237 VI. Geographical position and seasons (sh.n. ‘Lat-Long’ and ‘Incident light’): DS 15 is based on DS 5 (which  
 238 showed the most promising results) and includes geographical coordinates as additional drivers in the form  
 239 of sin(lat), sin(long), cos(long). DS 16 includes in addition to the DS 5 the role of incident light.
- 240 VII. Use of only PFTs and chlorophyll *a* (sh.n. ‘PFT only + CHL’): DS 17 is based on only the 12 PFTs, while  
 241 DS 18 is formed from DS 17 plus information on chlorophyll *a* averaged over the MLD. DSs 19 and 20 are  
 242 based on DS 6. To form the DS 19 we exclude temperature, NO<sub>3</sub> and PO<sub>4</sub> from the list of drivers in DS 6.  
 243 DS 20 is an extended version of DS 19 with all 12 PFTs concentration averaged over the MLD.
- 244 VIII. Chlorophyll *a* and chemical variables (sh.n. ‘Biochemical variables’): DSs 21, 22, 23, 24 are based on DS 5  
 245 and test the individual influence of chlorophyll *a* (DS 21), NO<sub>3</sub> (DS 22), PO<sub>4</sub> (DS 23) vertical profiles and  
 246 its ensemble (DS 24).
- 247 IX. Previous time step (sh.n. ‘Month - 1’): DSs, 25, 26 and 27 investigate the role of chlorophyll *a* (DS 27) and  
 248 some zooplanktons from the previous time step: gelatinous zooplankton and microzooplankton (DS26);  
 249 gelatinous zooplankton, micro- and macrozooplankton, averaged over MLD chlorophyll *a* and  
 250 coccolithophore (DS25).

251 The evaluation of the method is based on the mean correlation coefficient, total root-mean square errors (RMSE), and  
 252 total absolute bias between the ML outputs and PlankTOM12 POCs and POC<sub>L</sub> components. Moreover, we provide  
 253 the global maps of correlation coefficient and RMSE to vertical profiles of POC<sub>s</sub> and POC<sub>L</sub> at each grid point. Global  
 254 maps help to identify zones where the large errors can be hidden in the mean diagnostics due to the error compensation.

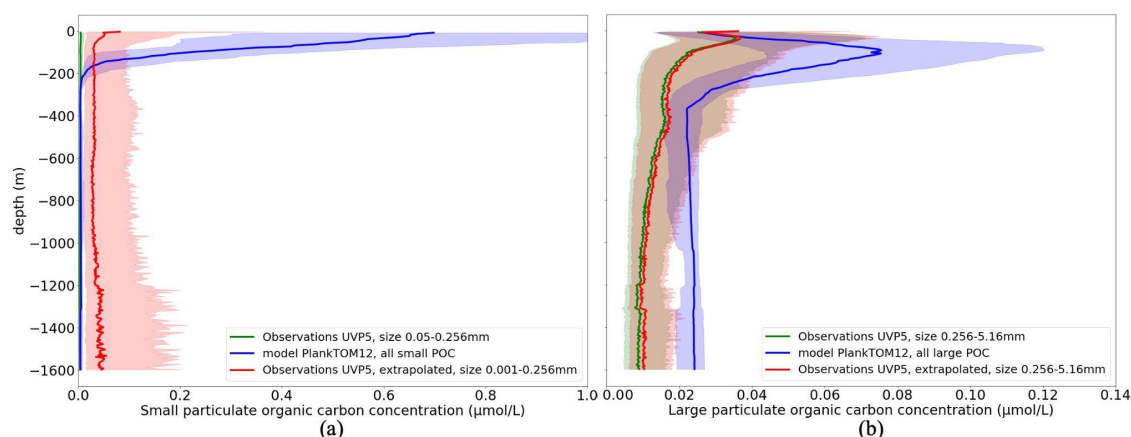
255 **3. Results.**  
 256



### 257 3.1. Data analysis.

258 In this study we test the capacity to reconstruct particulate organic carbon from sparse observations by using ML and  
259 a synthetic data set based on the PlankTOM12 model output. We compare observations and the output of the ocean  
260 model to provide a minimum of validation for the model data and to help explain differences in ML results when  
261 applied to real observations in the future.

262 Figure 4 shows the vertical profile of small ( $POC_S$ ) (Fig.4a) and large ( $POC_L$ ) particulate organic carbon (Fig.4b)  
263 based on the median from observations (green) and from daily PlankTOM12 model output (blue). Shading  
264 corresponds to values between 0.25 and 0.75 percentiles.



265 **Figure 4. Comparison of the vertical distribution of particulate organic carbon concentrations ( $\mu\text{mol/l}$ ) from UVP5**  
266 **measurements (green), PlankTOM12 daily model (blue) and extrapolated UVP5 measurements (red): (a) - small particulate**  
267 **organic carbon concentrations; (b) - large particulate organic carbon concentrations. The median is shown in dark and the**  
268 **shading corresponds to values between the 0.25 and 0.75 percentiles. The size of the particles does not correspond**  
269 **completely to the observations and the model, for  $POC_L$  the UVP particle range is chosen as 0.256-5.16 mm that**  
270 **corresponds approximately to the  $POC_L$  in the model.**  
271

272 PlankTOM12 overestimates  $POC_S$  up to  $3 \mu\text{mol/L}$  in the first 200m (Fig.4a, green and blue curves). UVP5 does not  
273 capture all small particles that is why we extrapolated the size range of UVP measurements (red curve, see details in  
274 2.1.1). The extrapolated measurements show an increase in  $POC_S$  in the first 100m, however this increase still results  
275 in the lower concentration compared with PlankTOM12. These results indicate that PlankTOM12 overestimates the  
276 concentration of small particulate organic carbon. PlankTOM12 also overestimates  $POC_L$  by up to  $0.08 \mu\text{mol/L}$   
277 in the first 200m and does not catch the increase in  $POC_L$  between 300 and 500m. Observations show an increase in  
278  $POC_L$  concentration in the first 50m while PlankTOM12 reproduces it lower, at 100m. The RMSE between modelled  
279 and observed  $POC_S$  is  $0.33 \mu\text{mol/L}$ , with correlation coefficient equals 0.083. RMSE equals  $0.23 \mu\text{mol/L}$  with  
280 correlation coefficient 0.061 for  $POC_L$ . The exclusion of isolated large values of  $POC_L$  ( $>2 \mu\text{mol/L}$ ) from the  
281 observation data set reduces the RMSE of  $POC_L$  to  $0.062 \mu\text{mol/L}$  with correlation 0.18. We believe that these  
282 differences result from differences in space and time resolution of observations and ocean model outputs. *In-situ*  
283 measurements are obtained at a particular time of the day and a particular latitude-longitude position while the model  
284 provides estimations over the day (or month) and on the model grid ( $2^\circ$  longitude and mean  $1.1^\circ$  latitude resolution).

285 We concluded that observed and modelled  $POC_S$  and  $POC_L$  have a common tendency in their vertical distributions.  
286 However, among other things, differences in amplitudes may affect our findings in this work when we develop a ML  
287 method based on observations only.

288 Due to the constraint in data availability further we use monthly PlankTOM12.

289 Before developing a ML method, we investigate the interactions between targets and drivers in the model. Table 2  
290 shows the correlation coefficients between the  $POC_S$  and  $POC_L$  and corresponding drivers that can influence  $POC_S$   
291 and  $POC_L$  variability.  $POC_S$  correlates with gelatinous zooplankton (GEL,  $r=0.66$ ), microzooplankton (PRO,  $r=0.63$ ),



292 coccolithophore (COC,  $r=0.56$ ), as well as with their values from previous time step (GEL,  $r=0.67$ ; PRO,  $r=0.51$ ;  
 293 COC,  $r=0.59$ ). Coccolithophore is one of the most abundant phytoplankton types in this version of the PlankTOM  
 294 model (similar to Wright et al., 2021). The growth of phytoplankton transfers dissolved inorganic carbon into dissolved  
 295 organic carbon which further aggregates into POC<sub>s</sub> and POC<sub>L</sub>. Also, POC<sub>s</sub> is generated from microzooplankton  
 296 egestion and excretion (Fig. 2). In addition to the mentioned above PFTs, POC<sub>s</sub> shows a correlation 0.44 with  
 297 temperature vertical profile at both the considered time step and at the previous time step. POC<sub>s</sub> has a negative  
 298 correlation with NO<sub>3</sub> ( $r=-0.46$ ) and PO<sub>4</sub> ( $r=-0.41$ ).

299 POC<sub>L</sub> does not show a high correlation with any of the proposed drivers individually and is therefore most likely the  
 300 result of multiple processes and/or multiple drivers, including for its production and destruction. The ML approach  
 301 should be able to identify combinations of drivers beyond straight correlations that are investigated directly here.  
 302 POC<sub>L</sub> has the highest correlation with chlorophyll *a* ( $r=0.42$ ), gelatinous zooplankton at the considered time step  
 303 ( $r=0.37$ ), and at previous time step ( $r=0.36$ ). Gelatinous zooplankton contribute to POC<sub>L</sub> formation through egestion  
 304 and excretion mainly from mucus (Fig. 2). As explained in Wright et al. (2021), mucus forms a large low-density mass  
 305 through aggregation with other particles. It can explain a correlation of gelatinous zooplankton with POC<sub>L</sub> in  
 306 PlankTOM12.

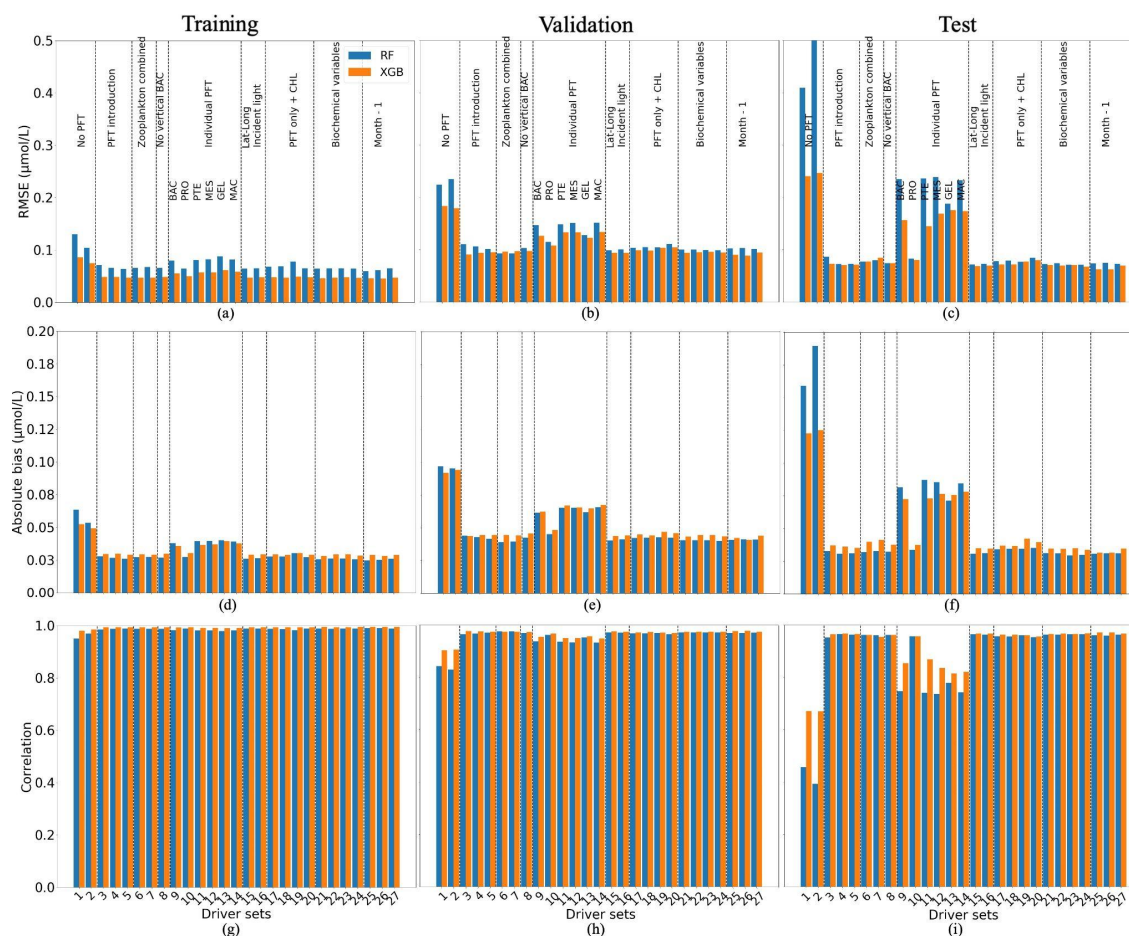
307 **Table 2. Correlation coefficient between small (POC<sub>s</sub>) and large (POC<sub>L</sub>) particulate organic carbon concentration and**  
 308 **possible drivers. Estimation is based on monthly PlankTOM12 output at the position of real-world observations from Fig.**  
 309 **1. ‘vp’ – vertical profile, ‘mean’ – average over MLD, ‘back’ – values from previous month.**  
 310

Driver	POC <sub>s</sub>	POC <sub>L</sub>	Driver	POC <sub>s</sub>	POC <sub>L</sub>	Driver	POC <sub>s</sub>	POC <sub>L</sub>
POC	1.00	0.33	BAC vp	-0.14	0.15	BAC back vp	-0.10	0.09
GOC	0.33	1.00	MES vp	-0.09	0.07	MES back vp	-0.09	-0.07
Depth	-0.32	-0.24	PTE vp	-0.07	0.17	PTE back vp	-0.08	0.08
Temperature vp	<b>0.44</b>	0.17	DIA vp	-0.04	0.15	DIA back vp	-0.03	0.09
Temp back vp	<b>0.44</b>	0.17	COC vp	<b>0.56</b>	0.31	COC back vp	<b>0.60</b>	0.31
MLD	-0.01	-0.07	PIC vp	0.00	0.07	PIC back vp	0.06	0.06
NO <sub>3</sub> vp	<b>-0.46</b>	0.01	PHA vp	0.27	0.15	PHA back vp	0.30	0.17
PO <sub>4</sub> vp	<b>-0.41</b>	0.04	GEL vp	<b>0.66</b>	<b>0.37</b>	GEL back vp	<b>0.68</b>	<b>0.36</b>
NO <sub>3</sub> back vp	<b>-0.46</b>	0.03	PRO vp	<b>0.63</b>	0.16	PRO back vp	<b>0.51</b>	0.14
PO <sub>4</sub> back vp	<b>-0.41</b>	0.05	MAC vp	0.07	0.14	MAC back vp	0.08	0.13
CHL vp	0.18	<b>0.42</b>	MIX vp	0.07	0.17	MIX back vp	0.03	0.05
CHL back vp	0.11	0.22	FIX vp	-0.00	0.23	FIX back vp	-0.00	0.23

311

312 **3.2. Development of the Machine Learning method.**





313

314 **Figure 5.** Comparison of the performance of the Random Forest (RF) and XGBoost methods and their fit to data for small  
 315 (POC<sub>s</sub>) particulate organic carbon concentration; (a, b, c) - RMSE in µmol/l, (d, e, f) - absolute bias in µmol/l, (g, h, i) -  
 316 correlation coefficient; (a, d, g) - training data set, (b, e, h) - the validation data set, (e, f, i) - the test data set. Results  
 317 compare data from the original (sampled) PlankTOM12 model output and POC<sub>s</sub> reconstructed using RF (blue) and XGB  
 318 (orange). The low RMSE and absolute biases indicate better performance of the ML method.

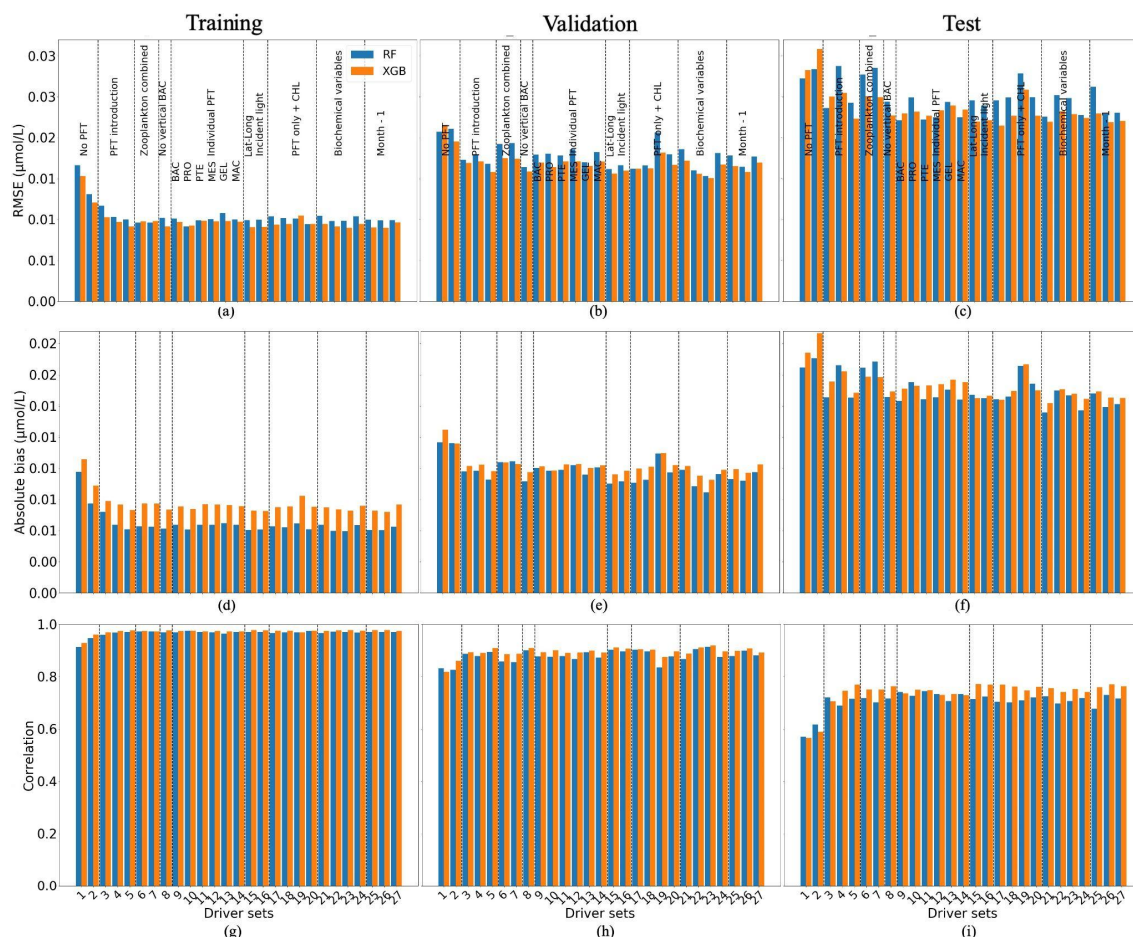
319 We tested 27 sets of drivers (Table 1) and two ML methods, Random Forest (RF) and XGBoost regression (XGB).

320 Figure 5 shows the statistics of POC<sub>s</sub> reconstruction using RF and XGB. XGB (orange) generally overperforms RF  
 321 (blue). The statistics are slightly worse for the validation and test data sets, as expected. For reconstructions using  
 322 XGB, the RMSE and absolute bias are about 0.05 µmol/L and 0.03 µmol/L on the training data set and vary around  
 323 0.1 µmol/L and 0.05 µmol/L, on the validation and test data, respectively. Correlation coefficients (Fig. 5g, h, i) have  
 324 high values on all datasets showing that the vertical profiles of POCs have a correct shape. These results show that  
 325 the available spatial and temporal coverage of *in situ* observations can be sufficient to reconstruct POCs with an  
 326 appropriate accuracy over the global ocean. The analysis of global maps (shown below) will help to identify areas  
 327 with low accuracy and their differences with training regions.

328 The worse results (highest RMSE, highest absolute bias, lowest correlation) are produced when there are no PFTs in  
 329 the driver set (DS1 and DS2; Figure 5): for XGBoost, RMSEs are 0.24 µmol/L, absolute biases equal to 0.12 µmol/L  
 330 with correlation coefficient 0.67 on the test data sets. Poor results are also obtained for DS9, 11, 12, 13 and 14: these  
 331 5 driver sets do not have any information on microzooplankton (PRO) and show high RMSEs and absolute biases,



332 around 0.16  $\mu\text{mol/L}$  and 0.074  $\mu\text{mol/L}$ , with low correlation, 0.83, compared with other driver sets which include  
 333 PRO. These results indicate that microzooplankton plays an important role in  $\text{POC}_s$  variability in the PlankTOM12  
 334 model.



335  
 336 **Figure 6. Comparison of the performance of the Random Forest (RF) and XGBoost methods and their fit to data for large**  
 337 **( $\text{POC}_L$ ) particulate organic carbon concentration; (a, b, c) - RMSE in  $\mu\text{mol/l}$ , (d, e, f) - absolute bias in  $\mu\text{mol/l}$ , (g, h, i) -**  
 338 **correlation coefficient; (a, d, g) - training data set, (b, e, h) - the validation data set, (e, f, i) - the test data set. Results**  
 339 **compare data from the original (sampled) PlankTOM12 model output and  $\text{POC}_L$  reconstructed using RF (blue) and XGB**  
 340 **(orange). The low RMSE and absolute biases indicate better performance of the ML method.**

341 Figure 6 shows the statistics of  $\text{POC}_L$  reconstruction using RF and XGB. XGBoost again slightly overperforms RF on  
 342 most driver sets. Results for driver sets with PFTs show lower RMSEs and absolute biases, and higher correlation  
 343 coefficients. Except for the effect of PFTs on the  $\text{POC}_L$  reconstruction, we did not observe a clear influence of one  
 344 driver or group of drivers. Using XGBoost the reconstruction of  $\text{POC}_L$  shows the RMSE in DS1 is high at 0.03  $\mu\text{mol/L}$ ,  
 345 while it is in the range of 0.021-0.026  $\mu\text{mol/L}$  in DS3-DS27, with absolute bias in DS1 of 0.02  $\mu\text{mol/L}$  and 0.015-  
 346 0.018  $\mu\text{mol/L}$  for DS3-DS27 based on test data (Fig. 6c, f). Likewise, a correlation coefficient of 0.56 for DS1, and  
 347 between 0.7 and 0.77 for DS3-DS27 based on the training data set (Fig. 6g).

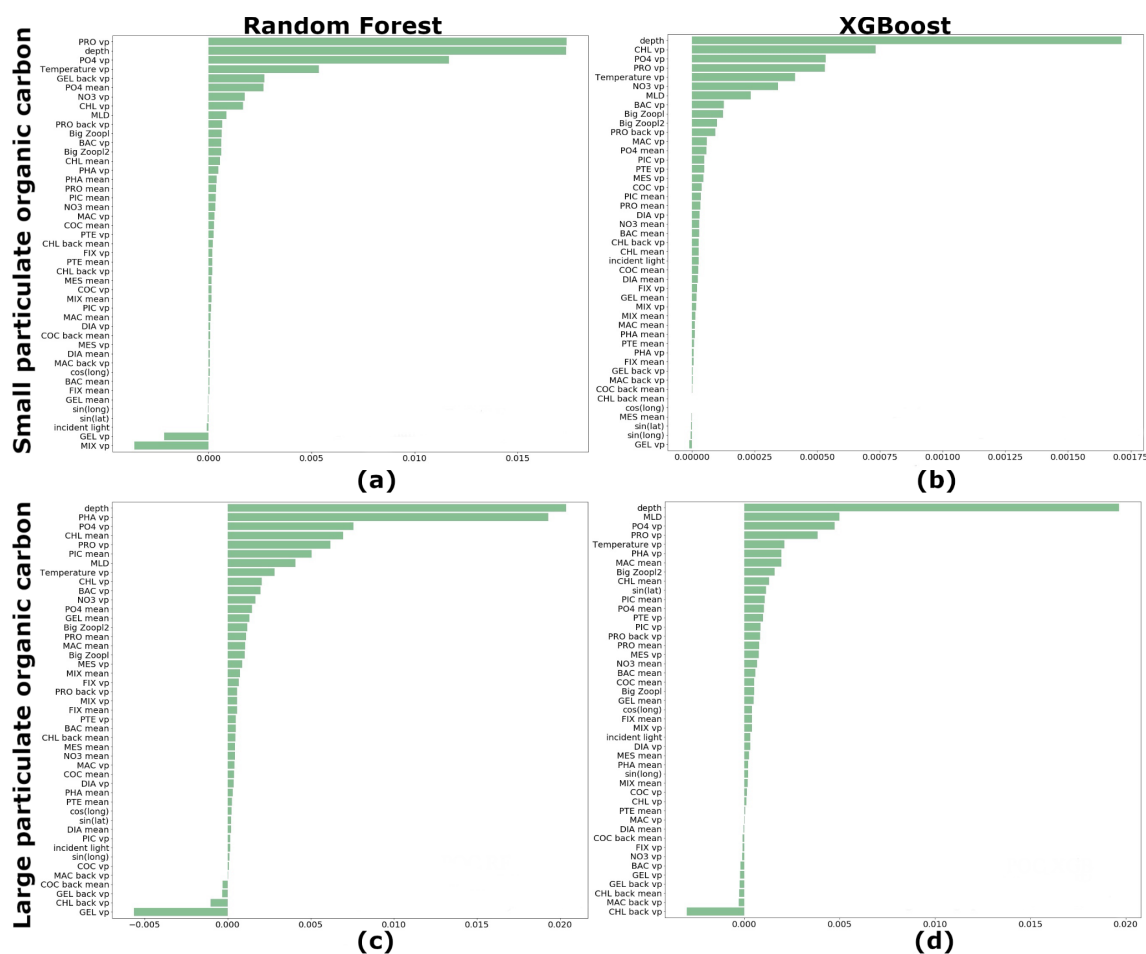
348 We estimated the ranking of importance for each driver averaged over 27 driver sets (Table 1) for RF and XGB (Fig.  
 349 7). Both, RF and XGB, show that microzooplankton (PRO), depth level, temperature,  $\text{NO}_3$  and  $\text{PO}_4$  play a dominant  
 350 role in reconstruction of  $\text{POC}_s$ . The absence of gelatinous zooplankton (GEL) can slightly improve the reconstruction.



351 Also, latitude and longitude do not affect POC<sub>S</sub> reconstruction. The depth level, temperature, MLD, microzooplankton  
 352 (PRO) and phaeocystis (PHA), PO<sub>4</sub>, and chlorophyll *a* averaged over MLD play a dominant role in POC<sub>L</sub>  
 353 reconstruction.

354 The sinus of latitude is in the top ten drivers that most affect POC<sub>L</sub> using XGBoost method: POC<sub>L</sub> distribution depends  
 355 on latitude zones. As for POC<sub>S</sub>, gelatinous zooplankton (GEL) shows a negative rank of driver importance and its  
 356 removal from the list of drivers can improve the statistics of reconstruction. Also, chlorophyll *a* concentration from  
 357 the previous month shows a similar effect on POC<sub>L</sub> (Fig. 7c, d).

358 Based on Figures 5, 6 and 7 we have chosen 10 driver sets with low RMSEs and absolute biases, and high correlation  
 359 coefficients (based on test data set) for POC<sub>S</sub> and POC<sub>L</sub> to provide global maps of these statistics and to see their  
 360 regional distributions. DS 5, 15, 16, 21, 22, 23, 24, 25, 26, 27 were chosen for further investigation of POC<sub>S</sub>  
 361 reconstruction; DS 5, 8, 15, 16, 17, 21, 23, 25, 26, 27 – for POC<sub>L</sub> reconstruction. Common for POC<sub>S</sub> and POC<sub>L</sub> driver  
 362 sets 5, 15, 16, 21, 23, 25, 26, 27 include all PFTs and their average over MLD, geographical positions and incident  
 363 light as well as chlorophyll *a*, PO<sub>4</sub>, and gelatinous zooplankton and microzooplankton from the previous time step  
 364 (Table 1). Also, we found that POC<sub>S</sub> reconstructions rest on biochemical conditions (DSs 21 and 24), while POC<sub>L</sub>  
 365 reconstruction mostly depends on the composition of the PFTs in the driver set (DSs 8 and 17). Additionally, we keep  
 366 DS1 to demonstrate a global effect of PFTs on reconstruction.



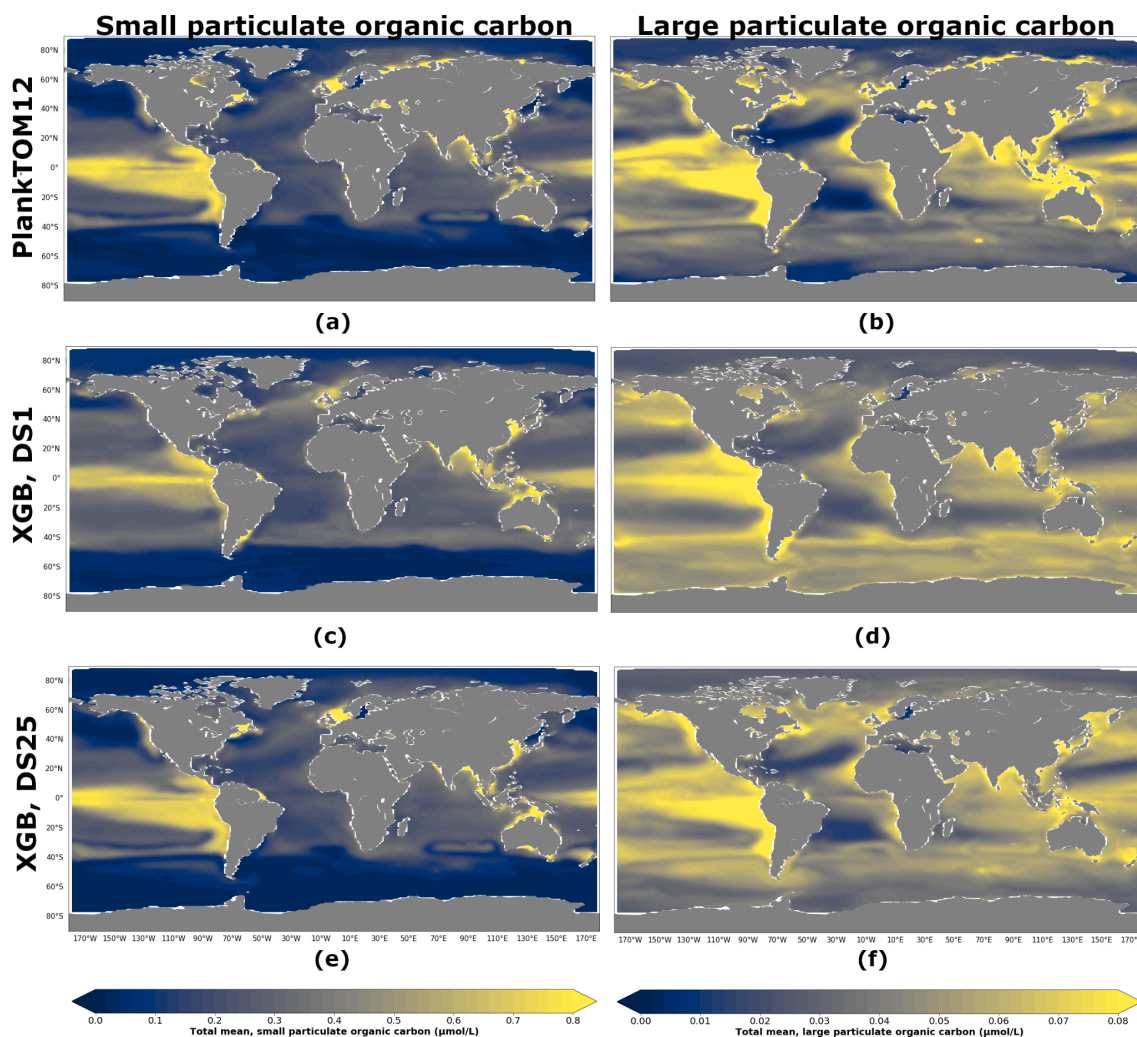
367  
 368 **Figure 7.** Ranking of importance for each driver averaged over 27 driver sets: (a) - Random Forest (RF) for reconstruction  
 369 of small (POC<sub>S</sub>) particulate organic carbon concentration; (b) - XGBoost (XGB) for small (POC<sub>S</sub>) particulate organic



370 carbon concentration; (c) - RF for  $POC_L$  concentration; (d) - XGB for  $POC_L$  concentration. 'vp' – vertical profile, 'mean'  
371 – average over MLD, 'back' – values from previous month.

### 372 3.3. $POC_S$ and $POC_L$ vertical profile reconstruction over the global ocean

373 In the previous section we showed that XGBoost provides the best results for the reconstructions of  $POC_S$  and  $POC_L$ .  
374 Further we use this ML method. Here we will discuss the regional results of DS1 without PFTs and 10 best driver sets  
375 chosen for each target separately.



376  
377 **Figure 8. Total averaged over the depth and period 2009-2013 small ( $POC_S$ ) and large ( $POC_L$ ) particulate organic carbon**  
378 **concentration: (a) – PlankTOM12  $POC_S$ , (b) – PlankTOM12  $POC_L$ , (c) – reconstruction of  $POC_S$  based on DS1 (NoPFT)**  
379 **using XGBoost, (d) – reconstruction of  $POC_L$  based on DS1 using XGBoos, (e) - reconstruction of  $POC_S$  based on DS25**  
380 **(vertical profiles of zooplanktons, and zooplankton and phytoplankton averaged over MLD) using XGBoost, (f) -**  
381 **reconstruction of  $POC_L$  based on DS25 using XGBoost.**

382 Figure 8 shows  $POC_S$  and  $POC_L$  concentration averaged over the depth and period 2009-2013 for PlankTOM12 (Fig.  
383 8a, b), XGboost reconstruction based on DS1 (Fig. 8c, d) and XGBoost reconstruction based on DS25 (Fig.8 e, f).



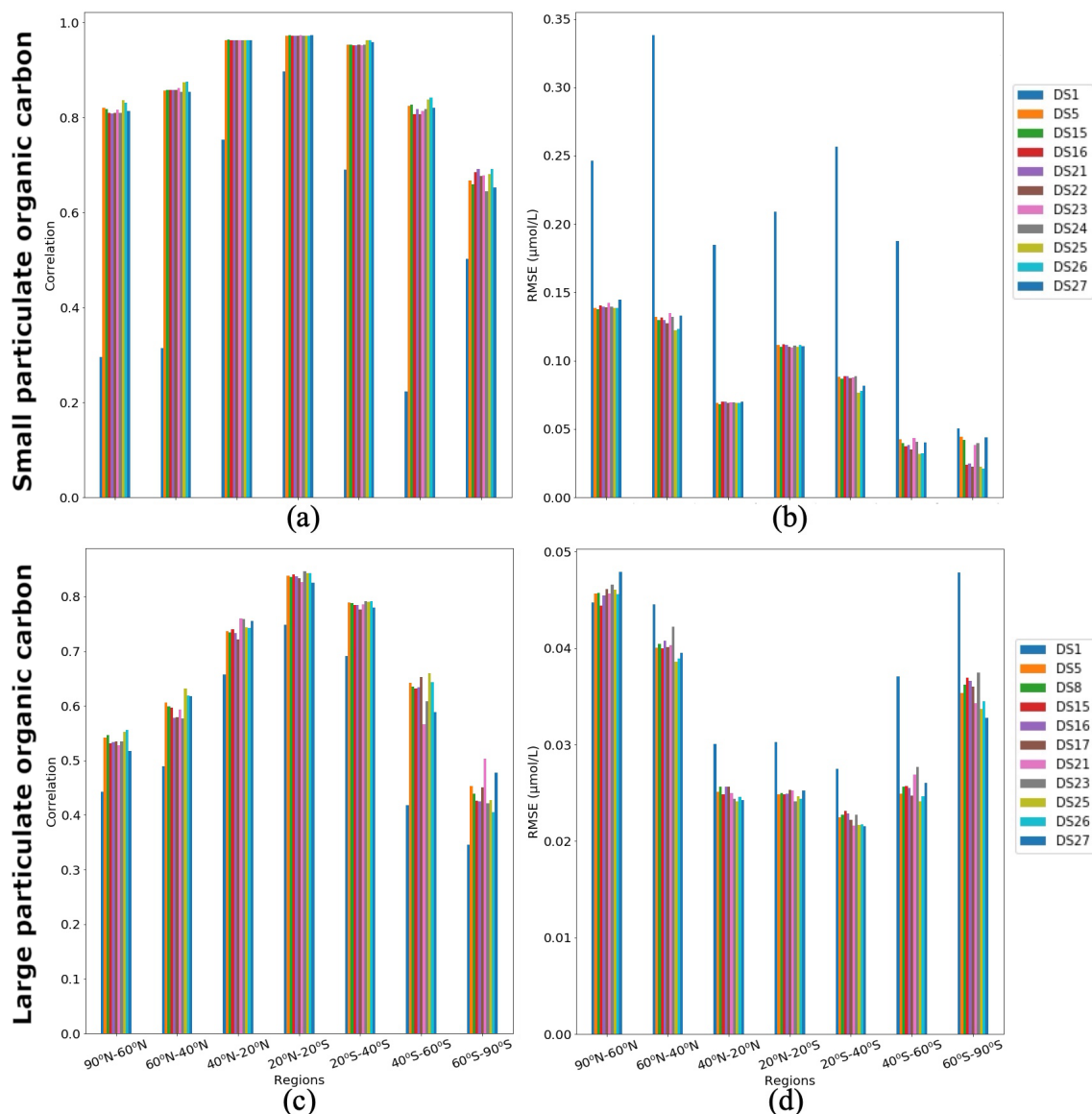


384 XGBoost captures well the spatial patterns: the high concentration of POCs in the Equatorial Eastern Pacific and its  
385 low concentration at high latitudes, as well as the high concentration of POC<sub>L</sub> in the Equatorial Eastern Pacific and in  
386 the North of the Indian Ocean and its low concentration in the Subtropical North and South Atlantic and in the  
387 Subtropical North Pacific. The presence of PFTs in driver sets (Fig. 8e, f) improves the reconstruction: the spatial  
388 patterns and its amplitude are visually close to ones from PlankTOM12 (Fig. 8a, b). The high concentration of POCs  
389 in the Equatorial Eastern Pacific is represented better using DS25 compared with DS1 where the concentration in the  
390 latitude band 0°S-20°S along the Peru is overestimated. Also, small decreases of POCs in the Subtropical North and  
391 South Atlantic are captured better when we use DS25. Similar for POCs, the high concentration in the Equatorial  
392 Eastern Pacific is represented better using DS25 compared with DS1 where the concentration misses the small  
393 decrease between 20°N and 0°N. Also, small decreases of POC<sub>L</sub> in the Subtropical North and South Atlantic as well  
394 as in the Subtropical North Pacific are pronounced better with DS25.

395 Figure 9 shows regional correlation coefficients and RMSEs between PlankTOM12 and XGBoost reconstruction over  
396 the global ocean for 2009-2013. We averaged correlation coefficient and RMSEs over 7 latitude zones: 90°N-60°N,  
397 60°N-40°N, 40°N-20°N, 20°N-20°S, 20°S-40°S, 40°S-60°S, 60°S-90°S. In POCs reconstruction, the DS1 shows the  
398 lowest correlation across latitude bands (between 0.22 and 0.9), and highest RMSEs (0.05-0.34 μmol/L; Fig.9a, b).  
399 DSs 25 and 26 show the highest correlations in the range of 0.68 (in region 60°S-90°S) and 0.97 (in region 20°N-20°S)  
400 and the lowest RMSEs in the range of 0.021 (in region 60°S-90°S) and 0.14 μmol/L (in region 90°N-60°N). DS25  
401 contains information on the previous-month distribution for micro-, macrozooplankton and gelatinous zooplankton  
402 vertical profiles as well as coccolithophores and chlorophyll *a* averaged over the MLD. DS26 is like DS25 but the  
403 drivers which bring information from the previous month are microzooplankton and gelatinous zooplankton vertical  
404 profiles.

405 10 driver sets (excluding DS1) show their highest RMSEs in POCs reconstruction in the region 90°N-60°N, with  
406 values up to 0.14 μmol/L in DS27 (Fig. 9b). Figure 10 shows maps of RMSEs (a, b) and correlation coefficients (c,  
407 d) between PlankTOM12 and reconstructed small particulate organic carbon (POCs) by XGBoost using driver sets 1  
408 (a, c) and 25 (b, d). The region 90°N-60°N shows improvement in RMSEs and absolute biases in DS25 compared with  
409 DS1, with RMSEs decreasing from 0.2 μmol/L to 0.03 μmol/L in Norwegian Sea, Baffin Bay, and the Arctic Ocean.  
410 However, errors stay high in the coastal regions, Northwestern passage and Hudson Bay that contribute to the high  
411 total RMSEs in this region. Results are similar for the region 60°N-40°N, where correlation coefficients increased  
412 from 0.3 to 0.87 on average over these zones (Fig. 10c, d). The tropical region 20°N-20°S shows correlation coefficient  
413 up to 0.97 for all driver sets except DS1. However, RMSEs are high in the tropical region, about 0.11 μmol/L on  
414 average (Fig. 9b), with RMSEs values of 0.2 μmol/L in the Tropical Eastern Pacific and Bay of Bengal in DS25 (Fig.  
415 10b). The high RMSEs in the Tropical Eastern Pacific can indicate insufficient data in a region of high interannual  
416 variability to correctly reconstruct POCs distribution. The region of the Southern Ocean (>60°S) shows the lowest  
417 correlation coefficients (in the range of 0.64-0.69) and RMSEs (in the range 0.023-0.044 μmol/L) for POCs (Fig. 9a,  
418 b). The inclusion of PFTs in the driver set significantly improves the RMSE in the region around 40°S for small  
419 (POCs) particulate organic carbon. The statistics are improved by about 75% in the region 40°S-60°S with RMSE  
420 decreasing from 0.18 (DS1) to 0.03 (DS25) and the correlation coefficient increasing from 0.22 (DS1) to 0.84 (DS25),  
421 on average (Fig. 9a, b; Fig. 10). The improvements in the Southern region are related to the role of zooplankton in the  
422 carbon flux in this area (Le Quééré et al., 2016; Wright et al., 2021).





423  
 424 **Figure 9. Correlations and RMSE averaged over latitude zones between PlankTOM12 and XGBoost reconstruction over**  
 425 **the global ocean for 2009-2013: (a, c) - correlation coefficient, (b, d) - RMSE in  $\mu\text{mol/l}$  (b, d); (a, b) - small particulate**  
 426 **organic carbon (POC<sub>s</sub>), (c, d) - large particulate organic carbon (POC<sub>L</sub>).**

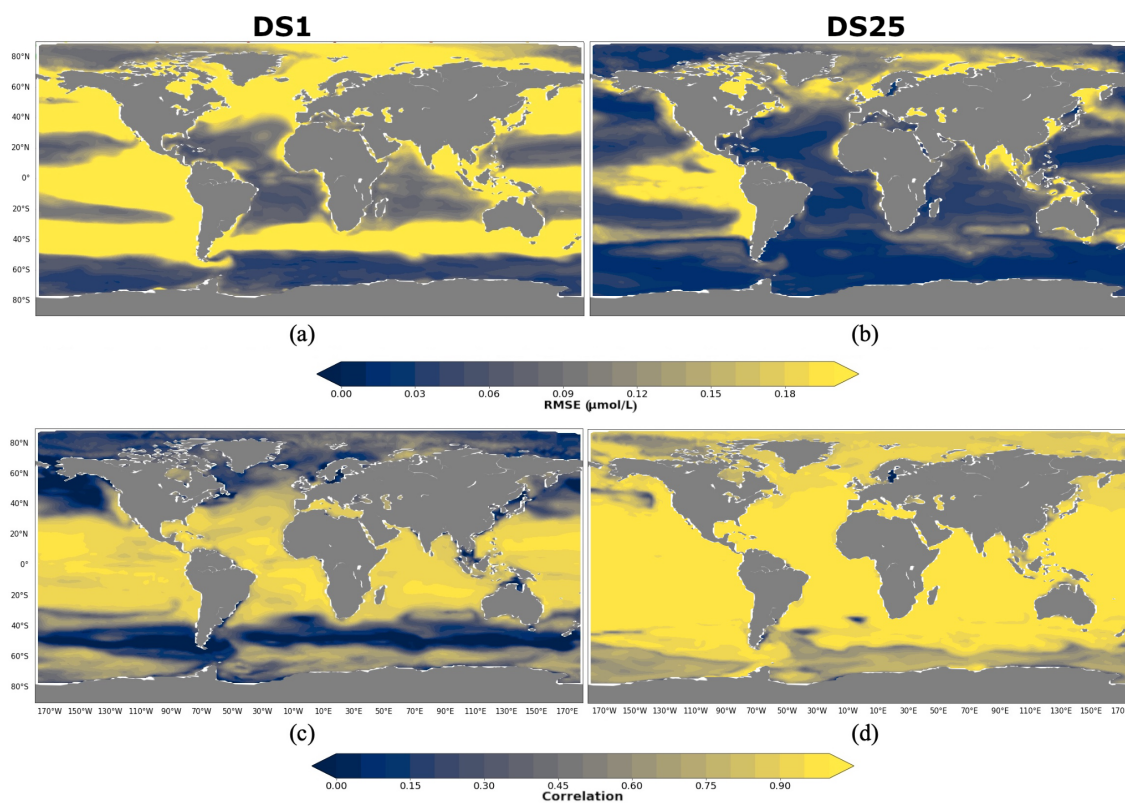
427 In POC<sub>L</sub> reconstruction, DS1 also shows the lowest correlation coefficients (0.35-0.75) and the highest RMSEs (0.027-  
 428 0.47  $\mu\text{mol/L}$ ) (Fig. 9c, d). DS25 shows the best results on average, with the correlation coefficient varying between  
 429 0.43 (in the region 60°S-90°S) and 0.84 (in the region 20°N-20°S), and RMSE varying between 0.021 (in the region  
 430 20°S-40°S) and 0.046 (in the region 90°N-60°N)  $\mu\text{mol/L}$ . POC<sub>L</sub> are reconstructed better in subtropical and tropical  
 431 regions compared to high latitude zones (Fig. 9c, d).

432 As for POC<sub>s</sub>, 10 driver sets (excluding DS1) show their highest RMSEs in POC<sub>L</sub> reconstruction in the region 90°N-  
 433 60°N, with values up to 0.05  $\mu\text{mol/L}$  in DS27 (Fig. 9d). Figure 11 shows maps of RMSEs (a, b) and correlation  
 434 coefficients (c, d) between PlankTOM12 and reconstructed large particulate organic carbon (POC<sub>L</sub>) by XGBoost using

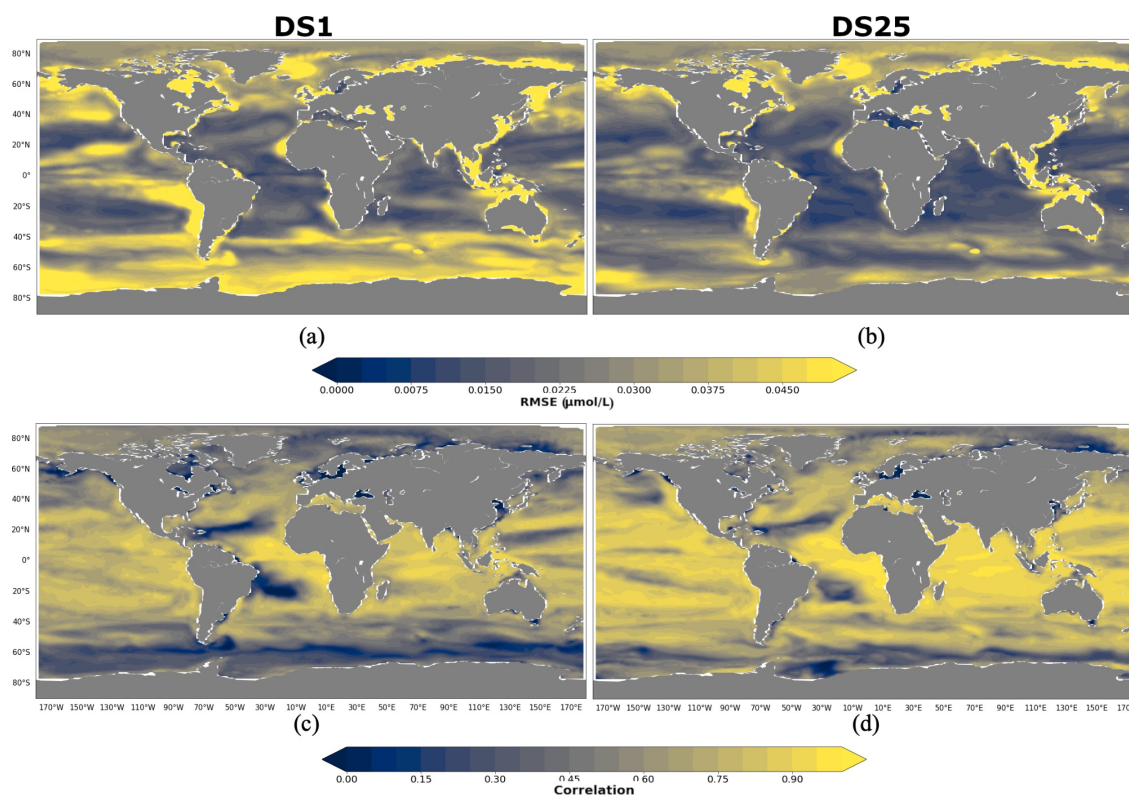


435 driver sets 1 (a, c) and 25 (b, d). Contrast to POC<sub>s</sub> reconstruction, the region 90°N-60°N does not show improvement  
436 in RMSEs for POC<sub>L</sub> reconstruction (Fig. 11b) in DS25 compared with DS1, with still high RMSEs in Norwegian Sea,  
437 Baffin Bay, and the Arctic Ocean, and additionally for POC<sub>L</sub> in Greenland Sea, where the algorithm did not have data  
438 for training. Similar to POC<sub>s</sub>, errors stay high in the coastal regions, Northwestern passage and Hudson Bay that  
439 contribute to the high total RMSEs in this region.

440 Global maps of statistics suggest that the most sensible region to driver set's composition for POC<sub>L</sub> is the Southern  
441 Ocean, as for POC<sub>s</sub> (Fig. 11). In the 40°S-60°S region, RMSE is reduced from 0.037 μmol/L in DS1 to 0.024 μmol/L  
442 in DS25 (Fig. 9d), and the correlation coefficient is increased from 0.42 to 0.66 (Fig. 9c) on average, respectively. In  
443 the Southern region 60°S-90°S, RMSE is reduced from 0.047 μmol/L in DS1 to 0.033 μmol/L in DS25, and the  
444 correlation coefficient is increased from 0.33 to 0.42 (Fig. 9c) on average, respectively. The average correlation  
445 coefficients in this zone were found to be less than 0.5 in all tests with the highest value 0.5 in DS21. DS21 contains  
446 all PFTs and chlorophyll *a* vertical profile as drivers. The RMSE for DS21 in this region is close to the one of DS25,  
447 0.34 μmol/L and 0.33 μmol/L, respectively. It identifies the importance of chlorophyll *a* in the Southern Ocean as  
448 driver of POC<sub>L</sub> variability.



449  
450 **Figure 10.** RMSE and correlation between monthly PlankTOM12 and results of POC<sub>s</sub> reconstruction using XGBoost over  
451 the period 2009-2013 for POC<sub>s</sub>. (a, b) – RMSEs, (c, d) – correlation coefficients; (a, c) – reconstruction based on DS1  
452 (NoPFT); (b, d) – reconstruction based on DS25 (vertical profiles of zooplanktons, and zooplankton and phytoplankton  
453 averaged over MLD).



454

455 **Figure 11. RMSE and correlation between monthly PlankTOM12 and results of  $\text{POC}_L$  reconstruction using XGBoost**  
456 **over the period 2009-2013 for  $\text{POC}_L$ . (a, b) – RMSEs, (c, d) – correlation coefficients; (a, c) – reconstruction based on**  
457 **DS1 (NoPFT); (b, d) – reconstruction based on DS25 (vertical profiles of zooplanktons, and zooplankton and**  
458 **phytoplankton averaged over MLD).**

459 The statistics of  $\text{POC}_S$  and  $\text{POC}_L$  reconstruction do not vary significantly between driver sets in all regions except in  
460 the Southern Ocean. This region is most sensitive to the composition of driver sets for both  $\text{POC}_S$  and  $\text{POC}_L$ .

#### 461 **4. Conclusion.**

462

463 The aim of this work was to test the potential of using Machine Learning to reproduce modelled concentrations of  
464 particulate organic carbon within the ocean using the distribution of available observations. We co-localised outputs  
465 of the PlankTOM12 global biogeochemical ocean model with the positions of observations of small ( $\text{POC}_S$ ) and large  
466 ( $\text{POC}_L$ ) particulate organic carbon concentrations. Using PlankTOM outputs as references we could identify the best  
467 ML method for POC reconstruction and estimate method's accuracy in regions with poor observational cover.

468 We tested two ML methods to reconstruct  $\text{POC}_S$  and  $\text{POC}_L$ : the XGBoost regressor and Random Forest. Both methods  
469 are algorithms based on decision trees. XGBoost overperformed Random Forest by about 9% on average for  $\text{POC}_S$   
470 reconstruction and by about 3% on average for  $\text{POC}_L$  reconstruction. XGBoost regressor builds the model sequentially  
471 improving it at each iterative step. At each iteration, XGBoost regressor analyses the prediction and gives more weight  
472 to the data where the fit is still wrong. It is a good tool for an unbalanced data set, like in our case where the data of  
473 particulate organic carbon concentration are sparse in time and space.

474 We tested the influence of a wide range of environmental and ecosystem drivers on  $\text{POC}_S$  and  $\text{POC}_L$  reconstruction.  
475 The introduction of Plankton Functional Types (PFTs) in the driver set greatly improves the fit and shows a linkage  
476 between surface ecosystem structure and particulate organic carbon distribution within the ocean interior. We



477 improved the accuracy of POC<sub>S</sub> reconstruction by 59% on RMSE, 63% on absolute bias and by 52% on correlation  
478 by introducing Plankton Functional Types (PFTs) in the driver sets (from the comparison of DS1 and DS25). The  
479 presence of PFTs in the driver sets also improved the accuracy of POC<sub>L</sub> reconstruction by 22% on RMSE, absolute  
480 bias and correlation (from the comparison of DS1 and DS25). POC<sub>S</sub> variability mostly depends on the depth level,  
481 vertical profiles of microzooplankton, temperature and PO<sub>4</sub>. POC<sub>L</sub> variability depends on the depth level, MLD,  
482 chlorophyll *a* averaged over MLD, vertical profiles of temperature, microzooplankton, phaeocystis and PO<sub>4</sub>.  
483 Additionally, we identified that chlorophyll *a* in driver sets improves the POC<sub>L</sub> reconstruction in the Southern Ocean.

484 Despite the good accuracy over the global ocean on average, the statistics are worse in the coastal regions and in the  
485 Tropical Eastern Pacific. The coastal regions suffer from the lack of data to represent the coastal dynamics. Therefore  
486 the ML reconstructions assign open-ocean processes to coastal regions, leading to significant biases. The Tropical  
487 Eastern Pacific is a region of strong interannual variability and the sparse measurements in time make it harder to  
488 capture this variability correctly. Other regions with poor coverage by observations - the Eastern Indian Ocean, the  
489 Western Pacific Ocean and the Southern Ocean - show the statistics of reconstruction comparable to one from regions  
490 with a good cover - regions in the Atlantic Ocean. However, we found that the Southern Ocean is a more sensible  
491 region to the driver set's composition. The observational data is particularly sparse in this region and our analysis  
492 suggests that identifying the drivers of importance based on real dataset will be difficult.

493 Here we showed that the XGBoost regressor and Random Forest are suitable for this problem and can reconstruct  
494 modelled POC<sub>S</sub> and POC<sub>L</sub> with appropriate accuracy. This is evidenced from the globally averaged correlation  
495 coefficient up to 0.88 for POC<sub>S</sub> and 0.68 for POC<sub>L</sub>, and the globally averaged RMSE up to 20 % (0.08 μmol/L) of  
496 standard deviation of PlankTOM12 POC<sub>S</sub>, and 65% (0.028 μmol/L) of standard deviation of PlankTOM12 POC<sub>L</sub>. ML  
497 outputs represent well the spatial patterns of POC<sub>S</sub> and POC<sub>L</sub> distribution. However, the validity of the approach on  
498 observations is dependent on the availability of co-located information on the drivers of importance. For some drivers  
499 this should be possible (e.g. environmental conditions and chlorophyll *a*), while for other drivers information is more  
500 sparse (e.g. the PFTs). Our analysis suggests that additional PFT observations would help provide broader insights  
501 into the distribution of POC in the ocean. The next step of this work is to apply ML to real data using methods from  
502 the present study.

503 This study provides insights on the drivers that may be responsible for POC<sub>S</sub> and POC<sub>L</sub> variability and regional  
504 dependencies. However, the dependencies are simply returning the outcome of complex ecosystem processes among  
505 the drivers as represented in the PlankTOM12 model. Although these processes are based on current understanding  
506 and a broad range of observations (Le Quéré et al., 2016; Wright et al., 2021; Buitenhuis et al., 2019), they remain  
507 results from a model output. Observations could reveal different drivers that are important for POC<sub>S</sub> and POC<sub>L</sub>.  
508 Depending on data availability and its time and space resolution, the final product based on observations should  
509 provide new insights on the drivers that govern particulate organic carbon concentration in the real ocean.

510 **Data and code availability.** PlankTOM12 data used within this study are available at  
511 <https://doi.org/10.5281/zenodo.7324781>. UVP5 data can be found at <https://doi.org/10.1594/PANGAEA.924375> (R.  
512 Kiko et al., 2021). Codes for data preparation, development of machine learning methods and tests of different driver  
513 sets as well as codes that provide figures shown in the article can be found at <https://doi.org/10.5281/zenodo.7326992>.

514 **Author contribution.** All authors contributed to the development of the methodology. ADS, CLQ, ETB designed the  
515 experiments, and ADS carried them out. ADS developed codes and performed the simulations. ADS prepared the  
516 paper with contributions from all coauthors.

517 **Acknowledgements.** The authors would like to thank Jean-Olivier Irisson for his contribution in the development of  
518 the methodology. ADS, ETB and CLQ acknowledge support from Royal Society (grant RP\R1\191063), and NERC  
519 Marine Frontiers project (grant NE/V011103/1) for CLQ. RK acknowledges support via a “Make Our Planet Great  
520 Again” grant from the French National Research Agency within the “Programme d’Investissements d’Avenir” (grant  
521 no. ANR-19-MPGA-0012) and by the Heisenberg program of the German Science Foundation under project number  
522 469175784.

523 **References**





- 524 Alldredge, A.: The carbon, nitrogen and mass content of marine snow as a function of aggregate size, *Deep Sea*  
525 *Research Part I: Oceanographic Research Papers*, 45, 529–541, [https://doi.org/10.1016/S0967-0637\(97\)00048-4](https://doi.org/10.1016/S0967-0637(97)00048-4),  
526 1998.
- 527 Batten, S. D., Abu-Alhaila, R., Chiba, S., Edwards, M., Graham, G., Jyothibabu, R., Kitchener, J. A., Koubbi, P.,  
528 McQuatters-Gollop, A., Muxagata, E., Ostle, C., Richardson, A. J., Robinson, K. V., Takahashi, K. T., Verheye, H.  
529 M., and Wilson, W.: A Global Plankton Diversity Monitoring Program, *Front. Mar. Sci.*, 6, 321,  
530 <https://doi.org/10.3389/fmars.2019.00321>, 2019.
- 531 Biau, G.: Analysis of Random Forest model, *Journal of Machine Learning Research*, 13(38), 1063–1095, 2012.
- 532 Buitenhuis, E. T., Hashioka, T., and Quéré, C. L.: Combined constraints on global ocean primary production using  
533 observations and models: OCEAN PRIMARY PRODUCTION, *Global Biogeochem. Cycles*, 27, 847–858,  
534 <https://doi.org/10.1002/gbc.20074>, 2013.
- 535 Buitenhuis, E. T., Le Quéré, C., Bednaršek, N., and Schiebel, R.: Large Contribution of Pteropods to Shallow CaCO<sub>3</sub>  
536 Export, *Global Biogeochem. Cycles*, 33, 458–468, <https://doi.org/10.1029/2018GB006110>, 2019.
- 537 Denvil-Sommer, A., Gehlen, M., Vrac, M., and Mejia, C.: LSCE-FFNN-v1: a two-step neural network model for the  
538 reconstruction of surface ocean pCO<sub>2</sub> over the global ocean, *Geosci. Model Dev.*, 12, 2091–2105,  
539 <https://doi.org/10.5194/gmd-12-2091-2019>, 2019.
- 540 Friedlingstein, P., Jones, M. W., O’Sullivan, M., Andrew, R. M., Bakker, D. C. E., Hauck, J., Le Quéré, C., Peters, G.  
541 P., Peters, W., Pongratz, J., Sitch, S., Canadell, J. G., Ciais, P., Jackson, R. B., Alin, S. R., Anthoni, P., Bates, N. R.,  
542 Becker, M., Bellouin, N., Bopp, L., Chau, T. T. T., Chevallier, F., Chini, L. P., Cronin, M., Currie, K. I., Decharme,  
543 B., Djeutchouang, L. M., Dou, X., Evans, W., Feely, R. A., Feng, L., Gasser, T., Gilfillan, D., Gkritzalis, T., Grassi,  
544 G., Gregor, L., Gruber, N., Gürses, Ö., Harris, I., Houghton, R. A., Hurtt, G. C., Iida, Y., Ilyina, T., Luijkx, I. T., Jain,  
545 A., Jones, S. D., Kato, E., Kennedy, D., Klein Goldewijk, K., Knauer, J., Korsbakken, J. I., Körtzinger, A.,  
546 Landschützer, P., Lauvset, S. K., Lefèvre, N., Lienert, S., Liu, J., Marland, G., McGuire, P. C., Melton, J. R., Munro,  
547 D. R., Nabel, J. E. M. S., Nakaoka, S.-I., Niwa, Y., Ono, T., Pierrot, D., Poulter, B., Rehder, G., Resplandy, L.,  
548 Robertson, E., Rödenbeck, C., Rosan, T. M., Schwinger, J., Schwingshackl, C., Séférian, R., Sutton, A. J., Sweeney,  
549 C., Tanhua, T., Tans, P. P., Tian, H., Tilbrook, B., Tubiello, F., van der Werf, G. R., Vuichard, N., Wada, C.,  
550 Wanninkhof, R., Watson, A. J., Willis, D., Wiltshire, A. J., Yuan, W., Yue, C., Yue, X., Zaehle, S., and Zeng, J.:  
551 Global Carbon Budget 2021, *Earth Syst. Sci. Data*, 14, 1917–2005, <https://doi.org/10.5194/essd-14-1917-2022>, 2022.
- 552 Friedrich, T. and Oschlies, A.: Basin-scale pCO<sub>2</sub> maps estimated from ARGO float data: A model study, *J. Geophys.*  
553 *Res.*, 114, C10012, <https://doi.org/10.1029/2009JC005322>, 2009.
- 554 Gorsky, G., Aldorf, C., Kage, M., Picheral, M., Garcia, Y., and Favole, J.: Vertical distribution of suspended  
555 aggregates determined by a new underwater video profiler, 1992.
- 556 Gorsky, G., Picheral, M., and Stemmann, L.: Use of the Underwater Video Profiler for the Study of Aggregate  
557 Dynamics in the North Mediterranean, *Estuarine, Coastal and Shelf Science*, 50, 121–128,  
558 <https://doi.org/10.1006/ecss.1999.0539>, 2000.
- 559 Guidi, L., Jackson, G. A., Stemmann, L., Miquel, J. C., Picheral, M., and Gorsky, G.: Relationship between particle  
560 size distribution and flux in the mesopelagic zone, *Deep Sea Research Part I: Oceanographic Research Papers*, 55,  
561 1364–1374, <https://doi.org/10.1016/j.dsr.2008.05.014>, 2008.
- 562 Guidi, L., Chaffron, S., Bittner, L., Eveillard, D., Larhlmi, A., Roux, S., Darzi, Y., Audic, S., Berline, L., Brum, J.  
563 R., Coelho, L. P., Espinoza, J. C. I., Malviya, S., Sunagawa, S., Dimier, C., Kandels-Lewis, S., Picheral, M., Poulain,  
564 J., Searson, S., Tara Oceans Consortium Coordinators, Stemmann, L., Not, F., Hingamp, P., Speich, S., Follows, M.,





- 565 Karp-Boss, L., Boss, E., Ogata, H., Pesant, S., Weissenbach, J., Wincker, P., Acinas, S. G., Bork, P., de Vargas, C.,  
566 Iudicone, D., Sullivan, M. B., Raes, J., Karsenti, E., Bowler, C., and Gorsky, G.: Plankton networks driving carbon  
567 export in the oligotrophic ocean, *Nature*, 532, 465–470, <https://doi.org/10.1038/nature16942>, 2016.
- 568 Hood, R. R., Laws, E. A., Armstrong, R. A., Bates, N. R., Brown, C. W., Carlson, C. A., Chai, F., Doney, S. C.,  
569 Falkowski, P. G., Feely, R. A., Friedrichs, M. A. M., Landry, M. R., Keith Moore, J., Nelson, D. M., Richardson, T.  
570 L., Salihoglu, B., Schartau, M., Toole, D. A., and Wiggert, J. D.: Pelagic functional group modeling: Progress,  
571 challenges and prospects, *Deep Sea Research Part II: Topical Studies in Oceanography*, 53, 459–512,  
572 <https://doi.org/10.1016/j.dsr2.2006.01.025>, 2006.
- 573 Kalnay, E., Kanamitsu, M., Kistler, R., Collins, W., Deaven, D., Gandin, L., Iredell, M., Saha, S., White, G., and  
574 Woollen, J.: The NCEP/NCAR 40-year reanalysis project, *B. Am. Meteorol. Soc.*, 77, 437–472, 1996.
- 575 Kiko, R., Picheral, M., Antoine, D., Babin, M., Berline, L., Biard, T., Boss, E., Brandt, P., Carlotti, F., Christiansen,  
576 S., Coppola, L., de la Cruz, L., Diamond-Riquier, E., de Madron, X. D., Elineau, A., Gorsky, G., Guidi, L., Hauss, H.,  
577 Irisson, J.-O., Karp-Boss, L., Karstensen, J., Kim, D., Lekanoff, R. M., Lombard, F., Lopes, R. M., Marec, C.,  
578 McDonnell, A., Niemeier, D., Noyon, M., O'Daly, S., Ohman, M. D., Pretty, J. L., Rogge, A., Searson, S., Shibata,  
579 M., Tanaka, Y., Tanhua, T., Taucher, J., Trudnowska, E., Turner, J. S., Waite, A. M., and Stemmann, L.: The global  
580 marine particle size distribution dataset obtained with the Underwater Vision Profiler 5 - version 1,  
581 <https://doi.org/10.1594/PANGAEA.924375>, 2021.
- 582 Kiko, R., Picheral, M., Antoine, D., Babin, M., Berline, L., Biard, T., Boss, E., Brandt, P., Carlotti, F., Christiansen,  
583 S., Coppola, L., de la Cruz, L., Diamond-Riquier, E., Durrieu de Madron, X., Elineau, A., Gorsky, G., Guidi, L.,  
584 Hauss, H., Irisson, J.-O., Karp-Boss, L., Karstensen, J., Kim, D., Lekanoff, R. M., Lombard, F., Lopes, R. M., Marec,  
585 C., McDonnell, A. M. P., Niemeier, D., Noyon, M., O'Daly, S. H., Ohman, M., Pretty, J. L., Rogge, A., Searson, S.,  
586 Shibata, M., Tanaka, Y., Tanhua, T., Taucher, J., Trudnowska, E., Turner, J. S., Waite, A., and Stemmann, L.: A  
587 global marine particle size distribution dataset obtained with the Underwater Vision Profiler 5, *Earth Syst. Sci. Data*  
588 *Discuss.* [preprint], <https://doi.org/10.5194/essd-2022-51>, 14, 4315–4337, <https://doi.org/10.5194/essd-14-4315-2022>, 2022.
- 590 Kirchman, D. L.: Growth Rates of Microbes in the Oceans, *Annu. Rev. Mar. Sci.*, 8, 285–309,  
591 <https://doi.org/10.1146/annurev-marine-122414-033938>, 2016.
- 592 Landschützer, P., Gruber, N., Bakker, D. C. E., Schuster, U., Nakaoka, S., Payne, M. R., Sasse, T. P., and Zeng, J.: A  
593 neural network-based estimate of the seasonal to inter-annual variability of the Atlantic Ocean carbon sink,  
594 *Biogeosciences*, 10, 7793–7815, <https://doi.org/10.5194/bg-10-7793-2013>, 2013.
- 595 Le Quéré, C., Harrison, S. P., Colin Prentice, I., Buitenhuis, E. T., Aumont, O., Bopp, L., Claustre, H., Cotrim Da  
596 Cunha, L., Geider, R., Giraud, X., Klaas, C., Kohfeld, K. E., Legendre, L., Manizza, M., Platt, T., Rivkin, R. B.,  
597 Sathyendranath, S., Uitz, J., Watson, A. J., and Wolf-Gladrow, D.: Ecosystem dynamics based on plankton functional  
598 types for global ocean biogeochemistry models, *Global Change Biology*, 0, 2016–2040,  
599 <https://doi.org/10.1111/j.1365-2486.2005.1004.x>, 2005.
- 600 Le Quéré, C., Buitenhuis, E. T., Moriarty, R., Alvain, S., Aumont, O., Bopp, L., Chollet, S., Enright, C., Franklin, D.  
601 J., Geider, R. J., Harrison, S. P., Hirst, A. G., Larsen, S., Legendre, L., Platt, T., Prentice, I. C., Rivkin, R. B., Saille, Y.,  
602 Sathyendranath, S., Stephens, N., Vogt, M., and Vallina, S. M.: Role of zooplankton dynamics for Southern Ocean  
603 phytoplankton biomass and global biogeochemical cycles, *Biogeosciences*, 13, 4111–4133,  
604 <https://doi.org/10.5194/bg-13-4111-2016>, 2016.
- 605 Lombard, F., Boss, E., Waite, A. M., Vogt, M., Uitz, J., Stemmann, L., Sosik, H. M., Schulz, J., Romagnan, J.-B.,  
606 Picheral, M., Pearlman, J., Ohman, M. D., Niehoff, B., Möller, K. O., Miloslavich, P., Lara-Lpez, A., Kudela, R.,  
607 Lopes, R. M., Kiko, R., Karp-Boss, L., Jaffe, J. S., Iversen, M. H., Irisson, J.-O., Fennel, K., Hauss, H., Guidi, L.,



- 608 Gorsky, G., Giering, S. L. C., Gaube, P., Gallager, S., Dubelaar, G., Cowen, R. K., Carlotti, F., Briseño-Avena, C.,  
609 Berlina, L., Benoit-Bird, K., Bax, N., Batten, S., Ayata, S. D., Artigas, L. F., and Appeltans, W.: Globally Consistent  
610 Quantitative Observations of Planktonic Ecosystems, *Front. Mar. Sci.*, 6, 196,  
611 <https://doi.org/10.3389/fmars.2019.00196>, 2019.
- 612 Mutshinda, C., Finkel, Z., Widdicombe, C., and Irwin, A.: Phytoplankton traits from long-term oceanographic time-  
613 series, *Mar. Ecol. Prog. Ser.*, 576, 11–25, <https://doi.org/10.3354/meps12220>, 2017.
- 614 Picheral, M., Guidi, L., Stemmann, L., Karl, D. M., Iddaoud, G., and Gorsky, G.: The Underwater Vision Profiler 5:  
615 An advanced instrument for high spatial resolution studies of particle size spectra and zooplankton: Underwater vision  
616 profiler, *Limnol. Oceanogr. Methods*, 8, 462–473, <https://doi.org/10.4319/lom.2010.8.462>, 2010.
- 617 Sauzède, R., Claustre, H., Uitz, J., Jamet, C., Dall’Olmo, G., D’Ortenzio, F., Gentili, B., Poteau, A., and Schmechtig,  
618 C.: A neural network-based method for merging ocean color and Argo data to extend surface bio-optical properties to  
619 depth: Retrieval of the particulate backscattering coefficient, *J. Geophys. Res. Oceans*, 121, 2552–2571,  
620 <https://doi.org/10.1002/2015JC011408>, 2016.
- 621 Sauzède, R., Bittig, H. C., Claustre, H., Pasqueron de Fommervault, O., Gattuso, J.-P., Legendre, L., and Johnson, K.  
622 S.: Estimates of Water-Column Nutrient Concentrations and Carbonate System Parameters in the Global Ocean: A  
623 Novel Approach Based on Neural Networks, *Front. Mar. Sci.*, 4, 128, <https://doi.org/10.3389/fmars.2017.00128>, 2017.
- 624 Sauzède, R., Johnson, J. E., Claustre, H., Camps-Valls, G., and Ruescas, A. B.: ESTIMATION OF OCEANIC  
625 PARTICULATE ORGANIC CARBON WITH MACHINE LEARNING, *ISPRS Ann. Photogramm. Remote Sens.*  
626 *Spatial Inf. Sci.*, V-2–2020, 949–956, <https://doi.org/10.5194/isprs-annals-V-2-2020-949-2020>, 2020.
- 627 Schlitzer, R.: Carbon export fluxes in the Southern Ocean: results from inverse modeling and comparison with  
628 satellite-based estimates, *Deep Sea Research Part II: Topical Studies in Oceanography*, 49, 1623–1644,  
629 [https://doi.org/10.1016/S0967-0645\(02\)00004-8](https://doi.org/10.1016/S0967-0645(02)00004-8), 2002.
- 630 Sunagawa, S., Acinas, S. G., Bork, P., Bowler, C., Tara Oceans Coordinators, Acinas, S. G., Babin, M., Bork, P.,  
631 Boss, E., Bowler, C., Cochrane, G., de Vargas, C., Follows, M., Gorsky, G., Grimsley, N., Guidi, L., Hingamp, P.,  
632 Iudicone, D., Jaillon, O., Kandels, S., Karp-Boss, L., Karsenti, E., Lescot, M., Not, F., Ogata, H., Pesant, S., Poulton,  
633 N., Raes, J., Sardet, C., Sieracki, M., Speich, S., Stemmann, L., Sullivan, M. B., Sunagawa, S., Wincker, P., Eveillard,  
634 D., Gorsky, G., Guidi, L., Iudicone, D., Karsenti, E., Lombard, F., Ogata, H., Pesant, S., Sullivan, M. B., Wincker, P.,  
635 and de Vargas, C.: Tara Oceans: towards global ocean ecosystems biology, *Nat Rev Microbiol*, 18, 428–445,  
636 <https://doi.org/10.1038/s41579-020-0364-5>, 2020.
- 637 Telszewski, M., Chazottes, A., Schuster, U., Watson, A. J., Moulin, C., Bakker, D. C. E., Gonzalez-Davila, M.,  
638 Johannessen, T., Kortzinger, A., Santana-Casiano, M., Wallace, D. W. R., and Wanninkhof, R.: Estimating the  
639 monthly pCO<sub>2</sub> distribution in the North Atlantic using a self-organizing neural network, 17, 2009.
- 640 Wright, R. M., Le Quéré, C., Buitenhuis, E., Pitois, S., and Gibbons, M. J.: Role of jellyfish in the plankton ecosystem  
641 revealed using a global ocean biogeochemical model, *Biogeosciences*, 18, 1291–1320, [https://doi.org/10.5194/bg-18-](https://doi.org/10.5194/bg-18-1291-2021)  
642 1291-2021, 2021.