# Global chlorophyll a concentrations of phytoplankton functional types with detailed uncertainty assessment using multi-sensor ocean color and sea surface temperature satellite products

Hongyan Xi<sup>1,1</sup>, Svetlana Loza (Losa)<sup>1</sup>, Antoine Mangin<sup>2</sup>, Philippe Garnesson<sup>2</sup>, Marine Bretagnon<sup>2</sup>, Julien Demaria<sup>2</sup>, Marinana A. Soppa<sup>1</sup>, Odile Hembise Fanton d'Andon<sup>2</sup>, and Astrid Bracher<sup>3</sup>

<sup>1</sup>Alfred Wegener Institute Helmholtz Centre for Polar and Marine Research <sup>2</sup>ACRI-ST <sup>3</sup>Alfred-Wegener-Institut für Polar und Meeresforschung, Universität Bremen, Universität Bremen

November 30, 2022

#### Abstract

Firstly, we re-tune an algorithm based on empirical orthogonal functions (EOF) for globally retrieving the chlorophyll *a* concentration (Chl-a) of phytoplankton functional types (PFTs) from multi-sensor merged ocean color (OC) products. The re-tuned algorithm, namely EOF-SST hybrid algorithm, is improved by: (i) using 30% more matchups between the updated global *in situ* pigment database and satellite remote sensing reflectance (Rrs) products, and (ii) including sea surface temperature (SST) as an additional input parameter. In addition to the Chl-a of the six PFTs (diatoms, haptophytes, dinoflagellates, green algae, prokaryotes and *Prochlorococcus*), the fractions of prokaryotes and *Prochlorococcus* Chl-a to total Chl-a (TChl-a), are also retrieved by the EOF-SST hybrid algorithm. Matchup data are further separated for low and high temperature regimes based on different PFT dependences on SST, to establish the SST-separated hybrid algorithms which further shows improved performance as compared to the EOF-SST hybrid algorithm. The per-pixel uncertainty of the retrieved TChl-a and PFT products is estimated by taking into account the uncertainties from both input data and model parameters through Monte Carlo simulations and analytical error propagation. The uncertainty assessment provided within this study sets the ground to extend the long-term continuous satellite observations of global PFT products by transferring the algorithm and its method to determine uncertainties to similar OC products until today. Satellite PFT uncertainty is also essential to evaluate and improve coupled ecosystem-ocean models which simulate PFTs, and furthermore can be used to directly improve these models via data assimilation.

- 1 Global chlorophyll *a* concentrations of phytoplankton functional types with
- 2 detailed uncertainty assessment using multi-sensor ocean color and sea
- 3 surface temperature satellite products
- Hongyan Xi<sup>1\*</sup>, Svetlana N. Losa<sup>1,2</sup>, Antoine Mangin<sup>3</sup>, Philippe Garnesson<sup>3</sup>, Marine
  Bretagnon<sup>3</sup>, Julien Demaria<sup>3</sup>, Mariana A. Soppa<sup>1</sup>, Odile Hembise Fanton d'Andon<sup>3</sup>, Astrid
  Bracher<sup>1,4</sup>
- <sup>1</sup> Alfred Wegener Institute, Helmholtz Centre for Polar and Marine Research, Bremerhaven,
   Germany
- 9 <sup>2</sup> Shirshov Institute of Oceanology, Russian Academy of Sciences, Moscow, Russia
- 10 <sup>3</sup> ACRI-ST, 06904 Sophia Antipolis Cedex, France
- <sup>4</sup> Institute of Environmental Physics, University of Bremen, Bremen, Germany
- 12 Correspondence to: H. Xi (Hongyan.Xi@awi.de)
- 13 Abstract:

14 Firstly, we re-tune an algorithm based on empirical orthogonal functions (EOF) for globally 15 retrieving the chlorophyll *a* concentration (Chl-a) of phytoplankton functional types (PFTs) 16 from multi-sensor merged ocean color (OC) products. The re-tuned algorithm, namely EOF-17 SST hybrid algorithm, is improved by: (i) using 30% more matchups between the updated 18 global *in situ* pigment database and satellite remote sensing reflectance (Rrs) products, and (ii) 19 including sea surface temperature (SST) as an additional input parameter. In addition to the 20 Chl-a of the six PFTs (diatoms, haptophytes, dinoflagellates, green algae, prokaryotes and 21 Prochlorococcus), the fractions of prokaryotes and Prochlorococcus Chl-a to total Chl-a 22 (TChl-a), are also retrieved by the EOF-SST hybrid algorithm. Matchup data are further 23 separated for low and high temperature regimes based on different PFT dependences on SST, 24 to establish the SST-separated hybrid algorithms which further shows improved performance 25 as compared to the EOF-SST hybrid algorithm. The per-pixel uncertainty of the retrieved 26 TChl-a and PFT products is estimated by taking into account the uncertainties from both input 27 data and model parameters through Monte Carlo simulations and analytical error propagation. 28 The uncertainty assessment provided within this study sets the ground to extend the long-term 29 continuous satellite observations of global PFT products by transferring the algorithm and its 30 method to determine uncertainties to similar OC products until today. Satellite PFT 31 uncertainty is also essential to evaluate and improve coupled ecosystem-ocean models which

32 simulate PFTs, and furthermore can be used to directly improve these models via data33 assimilation.

#### 34 Plain Language Summary:

35 Phytoplankton in the sunlit layer of the ocean contribute approximately 50% to global 36 primary production. They act as the base of the marine food web fueling fisheries, and also 37 regulate key biogeochemical processes such as exporting carbon to the deep ocean. 38 Phytoplankton contain various taxonomic groups that function differently in the marine 39 ecosystem. The global phytoplankton can be observed from space by analyzing the signal 40 leaving from the water surface recorded by the ocean color sensors onboard the satellites. 41 Based on an updated large global data set, satellite data from different ocean color sensors and 42 sea surface temperature data, we adapted our previous approach to better quantify the biomass 43 of the main six phytoplankton groups on the global scale. The uncertainty of the satellite 44 products of the phytoplankton groups are calculated by considering the errors propagated 45 from the satellite data and the model parameters. This approach for quantifying different 46 phytoplankton groups, together with the uncertainty assessment, can be extended to other 47 similar ocean color satellite data which cover different time periods, to ultimately generate 48 long term global distribution maps of multiple phytoplankton groups. This information will 49 help the modelers to better predict the phytoplankton changes in the future.

50 Keywords: Algorithm; empirical orthogonal functions; remote sensing reflectance; HPLC
51 pigments; merged products

#### 52 **1 Introduction**

53 Playing a fundamental role in the marine food web and biogeochemical cycling, 54 phytoplankton community structure and taxonomic composition has been widely investigated in the past decades, through various observations methods and ecological modelling (e.g., 55 56 Falkowski et al., 2003; Le Quéré et al. 2005; IOCCG, 2014). With a vast amount of quality-57 controlled ocean color remote sensing (OC) data, observations of the composition of 58 phytoplankton assemblages have provided various phytoplankton information including 59 dominance of phytoplankton groups, size classes (PSCs) and phytoplankton functional types 60 (PFTs) on a large scale both in space and time. The retrieval algorithms for phytoplankton composition were generally developed based on both in situ measurements and satellite 61 62 products, in which the former provides the ground truth information at specific time and 63 location, but with the latter the continuous spatiotemporal observations can be achieved.

Signals leaving the ocean surface (radiance or reflectance data) recorded by the satellites 64 65 inherit phytoplankton pigment information that relates to phytoplankton community structure and size classes (Bracher et al., 2017; Mouw et al., 2017). Therefore, they are often used to 66 67 establish spectral-based approaches to retrieve the concentrations of phytoplankton 68 chlorophyll, pigments and multiple PFTs from space (e.g., Alvain et al., 2005, 2008; Bracher 69 et al., 2009; Werdell et al., 2014; Correa-Ramirez et al., 2018; Xi et al., 2020; Lange et al., 70 2020). One of the efficient approaches is based on the empirical orthogonal function (EOF) 71 analysis on the spectral Rrs or water leaving radiance. By reducing the high dimensionality of 72 the spectral data, the dominant signals that best describe the variance of the structures lying in 73 the spectra can be assessed to establish the statistical models for predicting ocean color 74 metrics and various phytoplankton pigment and PFT chlorophyll a concentrations (Chl-a) 75 (e.g., Lubac and Loisel, 2007; Craig et al., 2012; Taylor et al., 2013; Bracher et al., 2015; 76 Soja-Woźniak et al., 2017; Xi et al., 2020; Lange et al., 2020). Approaches based on EOF 77 analysis also exhibited equivalent skill with little downgrading of the performance when 78 applied to reduced spectral resolution (Bracher *et al.*, 2015), enabling its wide applicability to 79 previous (e.g., SeaWiFS and MODIS) and current (e.g., MODIS, VIIRS and OLCI) 80 multispectral OC sensors and their merged products). In addition, as these approaches are 81 usually trained to retrieve ocean color metrics and PFT information directly from the satellite 82 spectral data, in its application no prior knowledge on the phytoplankton biomass or inherent 83 optical properties (IOPs) is required. This makes the implementation of such approaches 84 straightforward and practical for satellite OC products.

85 The EOF-based approach proposed by Xi et al. (2020) has justified to provide reliable Chl-a 86 retrievals of multiple PFTs on the global scale, through inter-comparisons with other satellite 87 derived PFT and PSC products. However, PFT retrievals by Xi et al. (2020) showed rather 88 low performance for prokaryotic phytoplankton. Incorporating additionally environmental 89 parameters, which are globally available from satellite measurements, such as optical depth, 90 sea surface temperature, wind stress, and light availability, have shown improvements to 91 several ocean colour PSC retrievals. For instance, Brewin et al. (2015) investigated the 92 influence of light in the mixed layer on the parameters of the three-component PSC model 93 (abundance-based model) of Brewin et al. (2010), and modified the model to better describe 94 the relationship between phytoplankton size structure and total chlorophyll with varying light 95 conditions. Ward (2015) and Brewin, Ciavatta, et al. (2017) both incorporated temperature 96 dependence into the three-component model and improved the model's ability in representing 97 Chl-a concentrations in all three PSCs using satellite estimates of SST and total Chl-a98 concentration.

99 There has been the emerging trend of the combined use of *in situ* data, satellite observations, 100 ecosystem modelling (Losa et al. 2019), as well as PFT or PSC data assimilation (Xiao and 101 Friedrichs, 2014; Pradhan et al., 2020), to allow comprehensive monitoring and predictions 102 of phytoplankton community structure. Satellite derived phytoplankton group-specific 103 products are also expected to be useful for validation of ecosystem model results (e.g., Ward 104 et al., 2012; Hirata et al., 2013; Holt et al., 2014; de Mora et al., 2016; Dutkiewicz et al. 2015; 105 Pradhan et al. 2019). One of the challenges to fulfill these tasks is associating the uncertainty 106 to the satellite derived PFT products (Bracher et al., 2017). Uncertainty estimates have been 107 well formulated and generated for other common OC algorithms that use satellite radiance 108 and reflectance data to derive OC products such as marine Chl-a concentration, diffuse 109 attenuation coefficient, and inherent optical properties (Werdell et al., 2018, McKinna et al., 110 2019). Though various approaches have been proposed to derive globally satellite 111 phytoplankton group products (Mouw et al. 2017), only the study by Brewin, Ciavatta, et al. 112 (2017a) has provided estimates of uncertainty on a per-pixel basis for the North Atlantic 113 Ocean. Uncertainty assessment can be carried out via two methods: validation through 114 comparison of the satellite retrievals with in situ data (e.g., Antoine et al. 2008; 115 Sathyendranath et al., 2019), or error propagation by accounting for the uncertainties in the 116 inputs and model parameters. Due to sparse distribution of the in situ measurements that 117 restrict the validation for the uncertainty estimation (Mélin and Franz, 2014), the error 118 propagation analysis has now been widely used to understand the sensitivity of the model 119 inputs and parameters to the outputs, and produce pixel-by-pixel uncertainty (e.g., Maritorena 120 et al., 2010; Lee et al., 2010; Kostadinov et al., 2016; Qi et al., 2017; Brewin, Tilstone, et al., 121 2017).

122 In this work, we improve the previously developed EOF-based algorithm of Xi et al. (2020) 123 for global retrievals of multiple PFT quantities by 1) including more matchup data between 124 the *in situ* pigment data set and satellite Rrs data from merged OC products, 2) accounting 125 for the sea surface temperature (SST) in the retrieval scheme, and 3) investigating the 126 influence of SST on the model parameters and the retrieved PFTs with the goal of 127 establishing a set of EOF-SST hybrid algorithms and improving the retrievals of TChl-a, Chl-128 a of six PFTs and the fractions of two prokaryotic phytoplankton. By applying the hybrid 129 algorithms to the merged OC products, we present the global distribution maps of the retrieved PFT quantities, and present a method to derive the per-pixel uncertainty propagated
from both the inputs and retrieval model for each satellite retrieved PFT quantity by
combining Monte Carlo (MC) simulations and an analytical approach.

### 133 **2 Data and Methods**

#### 134 **2.1 Data sets**

### 135 2.1.1 In situ data set of phytoplankton pigments

136 We updated the large global and open ocean (> 200 m) phytoplankton pigment data set 137 (spanning 2002-2012) from Losa et al. (2017) used in Xi et al. (2020) analyzed by High 138 Performance Liquid Chromatography (HPLC), by adding recently published HPLC pigment data (by February 2020) from SeaBASS, PANGAEA, British Oceanographic Data Centre 139 140 (BODC), and Open Access to Ocean Data (AODN) from Australia. All collected data were 141 quality controlled following the method by Aiken et al. (2009). A total of 9,595 sets of 142 pigment data were obtained as shown in Figure 1 with the distribution of total chlorophyll a concentration (TChl-a). In the database all required pigments for the PFT Chl-a calculation 143 144 were included (fucoxanthin, peridinin, 19'hexanoyloxy-fucoxanthin, 19'butanoyloxy-145 fucoxanthin, alloxanthin, total chlorophyll b, zeaxanthin and divinyl chlorophyll a).



Figure 1. Distribution of TChl-a (the sum of monovinyl chlorophyll *a*, divinyl chlorophyll *a*, chlorophyll *a* allomers, chlorophyll *a* epimers, and chlorophyllide *a*) from the quality controlled *in situ* pigment database (2002 - 2012) used in this study.

#### 150 2.1.2 Satellite data

#### 151 2.1.2.1 Satellite ocean color products

152 GlobColour data archive (http://www.globcolour.info/) has provided various OC products 153 from different sensors, including SeaWiFS, MODIS/AQUA, MERIS, VIIRS onboard Suomi-154 NPP and Sentinel-3A OLCI. In this study, we used the SeaWiFS-MODIS-MERIS merged 155 normalized remote sensing reflectance (Rrs) Level-3 (L3) product (hereafter referred to as 156 merged product) which covers the period from July 2002 to April 2012 from the GlobColour data archive (http://www.globcolour.info/; more details in ACRI-ST GlobColour Team et al., 157 158 2017). As in Xi et al. (2020), the daily merged product with 4-km resolution was used for 159 matchup extraction and monthly merged product with 25-km resolution was used for 160 algorithm application. Since this study focuses mainly on oceanic waters, shelf and coastal 161 waters (< 200 m) were masked out in the OC products using the ETOPO1 bathymetry 162 (Amante and Eakins, 2009).

163 2.1.2.2 Sea surface temperature (SST) data

164 The SST product used in this study was CMEMS OSTIA (Operational SST and Ice Analysis) 165 reprocessed analysis product, which is available on the CMEMS (Copernicus Marine 166 Environment Monitoring Service, https://marine.copernicus.eu/) platform, referenced as 167 SST\_GLO\_SST\_L4\_REP\_OBSERVATIONS\_010\_011. The CMEMS OSTIA reprocessed 168 analysis product is an interpolated product based on in situ measurements and satellite 169 observations from both infra-red and micro-wave data on a global regular grid at 0.05° 170 resolution (Donlon et al., 2012; Worsfold et al., 2020). The SST daily product from July 2002 171 to April 2012 was acquired and gridded to 4-km resolution. As for Rrs, monthly mean SST 172 product from 2002 to 2012 with 25-km resolution were also processed as input for deriving 173 global satellite PFT products.

#### 174 **2.1.3 Input data for PFT retrieval algorithm**

175 2.1.3.1 *In situ* PFT Chl-a concentration and fraction derived from diagnostic pigment analysis176 (DPA)

- 177 As described in Xi et al. (2020), Chl-a of PFTs were derived using an updated DPA method
- 178 (Soppa et al., 2014; Losa et al., 2017), that was originally developed by Vidussi et al., 2001,
- adapted in Uitz et al. (2006) and further refined by Hirata et al. (2011) and Brewin et al.
- 180 (2015). We used pigment concentrations from the *in situ* data base mentioned in Section 2.1.1
- 181 to derive the Chl-a of six PFTs diatoms, dinoflagellates, haptophytes, green algae,

182 prokaryotes, and Prochlorococcus. The partial coefficients of the diagnostic pigments used in 183 the DPA were the updated ones using a large global pigment data set as detailed in Table S1 184 in Losa et al. (2017), which were proved to be in good agreement with previous studies. Due 185 to the weak retrieval performance of prokaryotes and Prochlorococcus Chl-a in Xi et al. 186 (2020), in this study we included the fractions of prokaryotes (f-prokaryotes) and 187 Prochlorococcus (f-Prochlorococcus) to TChl-a, attempting to get improved retrievals of these two PFTs. PFT Chl-*a* lower than 0.005 mg m<sup>-3</sup> were excluded due to high uncertainty 188 189 (Xi et al., 2020) and the corresponding fractions of prokaryotes and Prochlorococcus were 190 also excluded.

191 2.1.3.2 Matchups between satellite SST and *in situ* PFT data

SST matchup data were extracted by matching spatially co-localized and temporally (on a daily basis) coincident with the *in situ* PFT measurements. A macro-pixel of 3×3 pixel centered on the *in situ* measurement was considered. If the standard deviation within this macro-pixel was lower than 25%, then the macro-pixel was considered suitable for the matchup. The median of the macro-pixel is defined as the SST value for the *in situ* site. Though the OSTIA SST product is quality controlled, abnormal values below -2 °C still existed, but were removed in the matchup data.

199 2.1.3.3 Matchups of Rrs merged product to *in situ* PFT and satellite SST data

200 Matchups of satellite R<sub>rs</sub> to in situ PFT data (which were also matchups to SST) were 201 extracted from global 4-km daily merged products. The same extraction and averaging 202 protocol including quality control as in Xi et al. (2020) were used to derive the single pixel 203 matchups. As justified in Xi et al. (2020), due to more matchup points and equivalent 204 retrieval performance compared to 3×3 matchups, the single pixel matchup data set was taken 205 as input data for the final retrieval approach (Figure 2). The Rrs matchup data with the nine 206 bands of 412, 443, 490, 510, 531, 547, 555, 670 and 678 nm from the merged products were chosen as the algorithm input data (Table 1). A total of 508 sets of matchup data covering the 207 208 global ocean were extracted (Figure 2).



Figure 2. Geographical locations of single pixel matchups of GlobColour merged  $R_{rs}$  at nine bands with *in situ* PFT and satellite SST data.

212 Table 1. List of the nine bands from sensors SeaWiFS, MODIS, and MERIS used in the

213 GlobColour merged products.

209

Sensors		Center Wavebands (nm)											
	412	443	490	510	531	547	555	670	678				
SeaWiFS	×	×	×	×			×	×					
MODIS	×	×	×		×	×	×	×	×				
MERIS	×	×	×	×			×*	×					

<sup>\*</sup> There was no band at 555 nm for MERIS itself, but the GlobColour Team provided for MERIS the 555

215 nm through an inter-spectral conversion made by using:

 $216 \qquad R_{rs}(555) = R_{rs}(560) * (1.02542 - 0.03757 * y - 0.00171 * y^2 + 0.0035 * y^3 + 0.00057 * y^4), \ where \ y = 0.00171 * y^2 + 0.0035 * y^3 + 0.00057 * y^4), \ y = 0.00171 * y^2 + 0.0035 * y^3 + 0.00057 * y^4), \ y = 0.00171 * y^2 + 0.0035 * y^3 + 0.00057 * y^4), \ y = 0.00171 * y^2 + 0.0035 * y^3 + 0.00057 * y^4), \ y = 0.00171 * y^2 + 0.0035 * y^3 + 0.00057 * y^4), \ y = 0.00171 * y^2 + 0.0035 * y^3 + 0.00057 * y^4), \ y = 0.00171 * y^2 + 0.0035 * y^3 + 0.00057 * y^4), \ y = 0.00171 * y^2 + 0.0035 * y^3 + 0.00057 * y^4), \ y = 0.00171 * y^2 + 0.0035 * y^3 + 0.00057 * y^4), \ y = 0.00171 * y^2 + 0.0035 * y^3 + 0.00057 * y^4), \ y = 0.00171 * y^2 + 0.0035 * y^3 + 0.00057 * y^4), \ y = 0.00171 * y^2 + 0.0035 * y^3 + 0.00057 * y^4), \ y = 0.00171 * y^2 + 0.0035 * y^3 + 0.00057 * y^4), \ y = 0.00171 * y^2 + 0.0035 * y^3 + 0.00057 * y^4), \ y = 0.00171 * y^2 + 0.00057 * y^4), \ y = 0.00$ 

217 log10(CHL1) and CHL1 is the total Chl-a concentration estimated by OC4 (ACRI-ST GlobColour Team et

218 *al.*, 2017). With this conversion, R<sub>rs</sub> at 555 nm for MERIS were also included in our study.

### 219 2.2 Algorithm re-tuning

220 The EOF-based PFT retrieval algorithm development and performance assessment were 221 detailed in Xi et al. (2020). Using the updated matchup data base, Xi et al. (2020) algorithm 222 was re-tuned and statistically assessed as detailed below. Figure 3 shows the scatterplots of 223 the matchup data of SST to TChl-a, the six PFT Chl-a and the two PFT fractions. Generally, 224 TChl-a, diatoms, haptophytes, dinoflagellates, and green algae show decreasing Chl-a with 225 increasing SST. However, Prokaryotes and Prochlorococcus Chl-a and their fractions to 226 TChl-a show positive correlation with SST. The statistically significant correlations indicate 227 that introducing SST as an additional term into the Xi et al. (2020) algorithm (see Section 228 2.2.1) might improve the algorithms' performance. Besides, at SST of around 8 °C, with

applying a 10-point running mean to the matchup data, there is a clear shift in the trends of 229 230 most PFT quantities as a function of SST. This further led us to establish for  $SST < 8 \ ^{\circ}C$  a different EOF-SST hybrid algorithm than for SST  $\geq 8$  °C (see Section 2.2.2). 231 232 *Prochlorococcus* data, as an exception from other PFTs, are rarely recorded in high latitudes 233 with low temperature (Flombaum *et al.*, 2013). Our matchup data set contained few (n = 6)234 divinyl chlorophyll a (a marker pigment of Prochlorococcus) measurements with low concentrations when SST < 8 °C. To construct the prediction models for all the PFTs more 235 easily, we excluded the regions where SST is below 8 °C for the Chl-a and fraction retrievals 236 237 of Prochlorococcus.



238

239 Figure 3. Scatterplots of in situ TChl-a, PFT Chl-a and fractions of prokaryotes and 240 Prochlorococcus versus collocated satellite SST data. The correlation coefficient R was 241 calculated based on the 10-point running mean (red curve).

242

#### 2.2.1 Adapted EOF-SST hybrid PFT algorithm based on the whole data set

243 The input data set used in the EOF-SST hybrid algorithm was the matchup data set that included the collocated nine-band Rrs from the merged products, SST satellite data and in situ 244 245 PFT data. Figure 4 depicts the flowchart of the EOF-SST hybrid algorithm, in which the EOF analysis remained unchanged by still using singular value decomposition (SVD) to decompose the (standardized) Rrs spectra into the EOF scores (U), singular values ( $\Lambda$ ) and EOF loadings (V) as in Xi *et al.* (2020). Now, when formulating the regression models of PFTs, we introduced SST as an additional term together with the column vectors  $u_{1,2,..,n}$  in U. Similar to Xi *et al.* (2020), we applied a stepwise routine to obtain the smaller regression model by removing least significant variables in U through minimization of the Akaike information criterion (AIC). The adapted regression model is expressed as

253 
$$\ln(C_p) = a_0 + a_1 u_1 + a_2 u_2 + \dots + a_{n} u_n + a_{SST} SST, \qquad (1)$$

254 where  $a_0$  is the intercept,  $a_{1,2,...n}$ , and  $a_{SST}$  are the regression coefficients for the selected EOF 255 scores and SST, respectively. With the adapted regression model, the same steps described as 256 cross validation and model assessment of Xi et al. (2020) were carried out to test the 257 robustness of the fitted model: the whole collocated data set was randomly split into two 258 subsets - the first subset containing 80% of the data was used for model fitting/training and 259 the rest 20% was used for prediction. The procedure was run for 500 permutations to 1) 260 record down in each permutation the model parameters for further uncertainty assessment, 261 and 2) generate a final statistical assessment based on the statistics of the model performance 262 derived from each permutation.

For the model assessment, we considered the slope (*S*), the intercept (*a*) of the GLM regression and coefficient of determination ( $\mathbb{R}^2$ ), which were based on the log-scaled PFT predictions against the log-scaled *in situ* PFT data. We also included the root-mean-square difference (RMSD), the median percent difference (MDPD), and the bias that were based on the non-log-transformed concentration data. These metrics were expressed as

269 MDPD = Median of 
$$\left[\frac{|(c_{pi}-c_{oi})|}{c_{oi}} \times 100\right], i = 1,...M,$$
 (3)

where *M* is the number of observations of PFTs ( $C_0$ ) and the corresponding predictions ( $C_p$ ). Meanwhile, the cross validation statistics ( $R^2cv$ , RMSDcv and MDPDcv), which represent both the model robustness and compromise model performance, was also determined by taking the mean of the statistical parameters  $R^2$ , RMSD and MDPD from all permutations, respectively.

#### 276 2.2.2 SST-separated hybrid PFT algorithms

Given the presented different SST-PFTs relationships for data set of SST < 8 °C and that of SST  $\geq$  8 °C (Figure 3), for all PFT quantities but *Prochlorococcus* Chl-a and fraction we separated the matchup data set based on 8 °C and established two specific EOF-SST hybrid algorithms using the two data sets (noted as SST-separated hybrid algorithms) following Section 2.2.1 (Figure 4). Note that the performance of the SST-separated hybrid algorithms was evaluated statistically based on all the predictions and *in situ* data to be consistent with that for the EOF-SST hybrid algorithm.

#### 284 **2.2.3** Application of algorithms

The established algorithms were applied to the satellite  $R_{rs}$  data from the merged products (Section 2.1.2) to retrieve PFTs globally (Figure 4). By projecting the  $R_{rs}$  data from the satellite onto the EOF loadings (V), a new set of EOF scores (U<sup>sat</sup>) was derived and was then used for the global PFT prediction together with the SST as an additional term in the fitted model in Eq. (1), where  $a_0$  and  $a_{1,2,...n}$  were obtained in the step of model training.



Figure 4. Flowchart illustrating the EOF-SST hybrid algorithm and the SST-separated hybrid algorithms for predicting TChl-a, Chl-a of six PFTs, and two fractions with GlobColour

293 merged product. The red dashed-line box depicts the model training with the pigment-satellite 294 matchup data; the green dashed-line box depicts the model application to satellite products 295 and the blue dashed-line box shows the output, i.e., the predicted PFT quantities.

#### 296 2.3 Uncertainty assessment of PFT retrieval

297 To quantify the uncertainty of the satellite PFT retrievals, we considered the uncertainties 298 propagated from the input datasets satellite  $R_{rs}$  ( $\sigma_{Rrs}$ ) and SST ( $\sigma_{SST}$ ), and uncertainty of the 299 model/algorithm parameters ( $\sigma_a$ ). Also other uncertainty sources exist, i.e., errors in the DPA 300 derived PFT data which results from the incorrect assignment of PFTs from marker pigments 301 and the *in situ* HPLC pigment measurement error. We were not able to obtain this information 302 for our large global data set collected from various cruises in the last decades. Mostly no other 303 descriptors of phytoplankton taxonomic composition had been measured and details on the 304 HPLC measurement error (including all associated steps, e.g., filtration, extraction and HPLC 305 analysis accuracy) are not available. According to IOCCG (2019), uncertainties of the HPLC-306 based chlorophyll a is around 7% and can be higher for other pigments (Claustre *et al.*, 2004). 307 In the current study, we could not quantify the combined uncertainty from both, HPLC 308 measurement and the DPA derived PFT, due to limited information therefore did not include 309 this error source in the uncertainty assessment.

310 All computations of the uncertainties in this study were based on the logarithmic transformed 311 data following conventional practice in the field of ocean color (OCCCI Product User Guide, 312 2020). However, we used natural logarithms instead of the common (base 10) logarithms, 313 because our algorithm was developed based on the natural logarithms. The common 314 logarithmic uncertainty can also be obtained by dividing our uncertainty by ln(10), i.e., 315 approximately a factor of 2.3. Due to the length of the article only the uncertainty derived 316 based on the whole EOF-SST hybrid algorithm is presented as a general approach to quantify 317 and consolidate the PFT uncertainty from different error sources.

#### 318 **2.3.1** Structure of the uncertainty propagation

319 With the EOF-SST hybrid retrieval models expressed in Eq. (2), the retrieval model applied 320 to the satellite data can be written in the following form:

321 
$$y(a, u(Rrs), SST) = \ln(C_P^{sat}) = a_o + a_1 u_1^{sat} + a_2 u_2^{sat} + \dots a_n u_n^{sat} + a_{SST} SST$$
 (5)

322 where *a* represents all the model coefficients; *u* represents all the EOF score vectors derived 323 from  $R_{rs}$  data therefore each can be expressed as:

$$u_i^{sat} = f(Rrs) \tag{6}$$

To estimate the final uncertainty of the retrieved PFTs,  $\sigma_y$ , we assume that the uncertainties due to *a*, *u*, and SST in Eq. (5) are not correlated with each other. According to the Guide to Uncertainty in Measurement (JCGM, 2008), the combined uncertainty of different sources could be estimated based on the law of propagation of uncertainty. Using the partial differences, the uncertainty of the PFT retrievals is presented theoretically as:

$$\sigma_{y} = \sqrt{\sigma_{y(Rrs)}^{2} + \sigma_{y(a)}^{2} + \sigma_{y(SST)}^{2}} = \sqrt{\sum_{i=1}^{N} \left(\frac{\partial y}{\partial Rrs_{i}}\right)^{2} \sigma_{Rrs_{i}}^{2} + \sum_{i=0}^{n} \left(\frac{\partial y}{\partial a_{i}}\right)^{2} \sigma_{a_{i}}^{2} + \left(\frac{\partial y}{\partial SST}\right)^{2} \sigma_{SST}^{2}}.$$

$$(7)$$

Since the uncertainties propagated from errors of model parameters ( $\sigma_{y(a)}$ ) and SST ( $\sigma_{y(SST)}$ ) are both linear, they can be analytically derived and expressed together as:

335 
$$\sigma_{y(a)}^2 + \sigma_{y(SST)}^2 = \sum_{i=1}^N (u_i^{sat})^2 \sigma_{a_i}^2 + SST^2 \sigma_{a_{SST}}^2 + (a_{sst})^2 \sigma_{SST}^2.$$
(8)

Where  $\sigma_{SST} = 0.46$  °C (Worsfold *et al.*, 2020 for the OSTIA SST product);  $\sigma_{a_i}$  and  $\sigma_{a_{SST}}$ 336 were determined during the cross validation procedure as described in Section 2.2.1. To 337 further understand how the  $\sigma_{a_i}$  and  $\sigma_{a_{SST}}$  were determined, Figure 5 shows the distributions of 338 339 the coefficients derived from all the permutations in the cross validation as an example for 340 diatoms. In general, the coefficient distributions followed the normal distribution, the 341 uncertainty of each coefficient was thus determined by calculating the corresponding standard 342 deviation. The uncertainties of the model coefficients for all other PFT quantities were also 343 determined in the same manner.

Since the converted prokaryotes and *Prochlorococcus* Chl-a were calculated by multiplying their fractions and the TChl-a together, the corresponding uncertainties were determined by the uncertainties of TChl-a and that of *f*-Prokaryotes (*f*-*Prochlorococcus*). Using  $y_{\text{proka}\_\text{conv}}$ ,  $y_1$ and  $y_2$  to denote the converted prokaryotes Chl-a, TChl-a and f-Prokaryotes, the uncertainty of the converted prokaryotes Chl-a,  $\sigma_{y_{\text{proka}\_\text{conv}}}$ , can be determined by the uncertainty of TChl-a and *f*-Prokaryotes, namely  $\sigma_{y_1}$  and  $\sigma_{y_2}$ , through the following equation:

350 
$$\sigma_{y_{proka\_conv}} = \sqrt{\sigma_{y_1}^2 + \sigma_{y_2}^2 + 2\sigma_{y_1}\sigma_{y_2}r_{12}},$$
 (9)

351 where  $r_{12}$  is the correlation coefficient between  $\sigma_{y_1}$  and  $\sigma_{y_2}$  (as both depend on SST-related 352 uncertainties). Eq. (9) also applies for the uncertainty of converted *Prochlorococcus* Chl-a.



353

Figure 5. Histograms depicting the distributions of the model coefficients derived from 500 permutations of the cross validation for diatoms.  $\sigma_{a_i}$  (i.e., standard deviation, SD) of each coefficient was determined accordingly.

Among all uncertainty components in Eq. (7), the uncertainty of the PFTs propagated from the errors of the satellite Rrs spectra,  $\sigma_{y(Rrs)}$ , is the challenging part to quantify, as it is nonlinear and not as straightforward as the other two uncertainty sources, due to the EOF analysis performed with the spectra. We therefore used a MC simulation-based approach to estimate the  $\sigma_{y(Rrs)}$  and detailed it in the following subsection 2.3.2.

### 362 2.3.2 Rrs uncertainty propagation

Based on the uncertainty of the water leaving radiance for SeaWiFS, MODIS-Aqua, and MERIS reported in Maritorena *et al.* (2010), Rrs absolute uncertainties for these sensors were derived and were used in GlobColour program. In our study, we took the root mean square (RMS) of the common bands from two or three sensors as the uncertainty of the merged products (Table 2). Using the matchup data for the merged products at nine bands, the following steps were carried out to fulfil the uncertainty propagation from Rrs to the PFTs.

369 1. The 508 Rrs matchup spectra were randomly divided equally into two datasets - 50%
370 as the training data set, and the other 50% as the testing data set (in a total of 254).

- 371 The corresponding matchups of *in situ* PFT and retrieved PFT data were also divided372 accordingly.
- 373 2. For the training data sets, we performed 10,000 MC simulations to randomly 374 introduce for each band the Rrs uncertainty ( $\sigma_{Rrs}$ ) to each Rrs spectrum in the training 375 data set (in a total of 2,540,000 simulated spectra).
- 376 3. The MC simulated Rrs spectra were applied to the EOF-SST hybrid algorithm to 377 estimate the PFTs with Rrs uncertainty taken into account. For each sample, 10,000 378 estimates of the PFT were generated from the 10,000 MC simulated Rrs, so that the 379 uncertainty (standard deviation,  $\sigma_{MC-PFT}$ ) of the PFTs were determined based on these 380 10,000 estimates for each specific sample.
- 381 4. When the  $\sigma_{MC}$  for all samples in the training data set were determined through Step 3, 382 a look-up table (LUT) was built for each PFT by fitting  $\sigma_{MC-PFT}$  as a function of the 383 retrieved PFT.
- 5. The LUT for each PFT was applied to the testing data set for the uncertainty validation, and also to the satellite PFT products to derive per-pixel uncertainty of the satellite PFT due to Rrs,  $\sigma_{y(Rrs)}$ , which was combined with uncertainties from the other sources via Eq. (7) to derive the final uncertainty of PFT satellite retrievals.
- Table 2. Absolute uncertainties of Rrs ( $\sigma_{\text{Rrs}}$ , Sr<sup>-1</sup>)) for different sensors in the merged products derived based on Maritorena *et al.* (2010). The root mean square (RMS) was taken as the uncertainty of the merged products.

Wavebands (nm)	) 412	443	490	510	531	547	555	670	678
MODIS $\sigma_{Rrs}$	0.00071	0.00063	0.00049	-	0.00024	0.00019	-	0.000055	0.000030
MERIS $\sigma_{Rrs}$	0.00066	0.00059	0.00047	0.00033	-	-	0.00023	0.00010	0.000098
SeaWiFS $\sigma_{Rrs}$	0.00072	0.00064	0.00050	0.00036	-	-	0.00025	0.000075	-
RMS $\sigma_{Rrs}$	0.00070	0.00062	0.00049	0.00035	0.00024	0.00019	0.00024	0.000080	0.000072

#### 391 2.3.3 Assessment of the per-pixel PFT uncertainty

With the steps in Section 2.3.2, the uncertainty propagated from the Rrs to the satellite retrieved PFTs was determined by applying the LUT to each pixel of the satellite derived PFT products, the term  $\sigma_{y(Rrs)}^2$  in Eq. (7) was thus derived. Together with the other two terms,  $\sigma_{y(a)}^2$  and  $\sigma_{y(SST)}^2$  that could be calculated analytically through Eq. (8), the combined PFT uncertainty  $\sigma_y$  of each pixel from different sources was ultimately obtained.

#### **397 3 Results and discussion**

### 398 **3.1 EOF-SST hybrid algorithms for PFT retrievals**

#### 399 **3.1.1** EOF-SST hybrid algorithm based on the whole matchup data set

400 Before setting up the EOF-SST algorithm, we firstly used the updated input data set to update 401 the original EOF-based algorithm proposed by Xi et al. (2020) where SST was not included, 402 the performance of the updated algorithm showed nearly identical performance as compared 403 to the original one presented in Xi et al. (2020) (details not shown). This indicates that the 404 original algorithm can be hardly improved by purely enlarging the training data set. Within the frame of the EOF-SST hybrid algorithm in the current study, TChl-a, PFT Chl-a and the 405 406 fractions of two PFTs were predicted based on the regression models built using the EOF 407 scores derived from the nine-band Rrs data, SST and the in situ PFT data. As presented in 408 Table 3 and Figure 6 (A-F, H), compared to the original algorithm in Xi et al. (2020), the 409 EOF-SST hybrid algorithm shows significant improvements for all predicted quantities 410 except for Prochlorococcus where weak performance still remains. For TChl-a, and Chl-a for diatoms, haptophytes, dinoflagellates and green algae,  $R^2$  and  $R^2$ cv are increased to 0.59 – 411 412 0.85 (compared to 0.51 - 0.76) and to 0.56 - 0.84 (compared to 0.47 - 0.75), respectively. 413 MDPD and MDPDcv are remarkably reduced to 30% to 55% and 31% to 56%, respectively, 414 for all quantities as compared to Xi et al. (2020) ranging 37%-74% and 37%-75%, 415 respectively. RMSD and RMSDcv values are also decreased significantly in the EOF-SST 416 hybrid algorithm compared to the previous results.

417 To further improve the prediction of prokaryotes and Prochlorococcus Chl-a, the hybrid 418 algorithm was also trained to retrieve the fractions of prokaryotes and Prochlorococcus to 419 TChl-a. This was motivated because prokaryotes dominate the low TChl-a mid- to low 420 latitude waters, so generally their Chl-a is low. By using their fraction instead of Chl-a a 421 better spread of the data is achieved that enhances the signal to be retrieved which is 422 beneficial for application in abundance-based PSC retrievals (e.g., Brewin et al., 2010). As expected, the prediction models for the two fractions performed well with  $R^2 > 0.62$  and 423 424 MDPD within 42% (Table 3). Though the overall performance of the fraction retrieval for the 425 two PFTs had been improved, the regressions between the predicted and observed fractions 426 (Figure 6G and J) show that higher discrepancies exist in low fraction values, indicating that 427 it is still difficult to deal with low Chl-a values. Using predicted TChl-a, the fractions were 428 further converted to Chl-a for the two PFTs. Table 3 shows that the converted prokaryotes

Chl-a retrieval displays much improved performance compared to the directly retrieved 429 prokaryotes (R<sup>2</sup> of 0.34 vs. 0.27, MDPD of 39% vs. 44%, and RMSD of 0.08 vs. 0.09 mg m<sup>-</sup> 430 <sup>3</sup>), but it is downgraded compared to the f-prokaryotes (Figure 6H). The fraction to Chl-a 431 432 conversion scheme shows barely improvement in predicting Prochlorococcus Chl-a (Table 3 433 and Figure 6K versus Figure 6I). Though *f-Prochlorococcus* is overall better predicted 434 compared to the direct retrieval of Prochlorococcus Chl-a, by using the conversion, the low signal to noise ratio in the retrieved TChl-a and *f*-Prochlorococcus deteriorates for  $R^2$  and 435 RMSD the final Prochlorococcus Chl-a estimation. The MDPD and bias are however slightly 436 437 improved. Weak prediction performance of the converted Prochlorococcus Chl-a reveals that 438 it is still challenging to enhance their retrieval accuracy to the same level as other PFTs due to

439 the low concentrations and small variability (Xi *et al.*, 2020).

Table 3. Statistics of regression models for TChl-a, six PFT Chl-a, fractions of prokaryotes 440 441 and Prochlorococcus and the corresponding converted Chl-a using SST and EOF modes 442 based on the nine-band Rrs matchups from merged OC products (upper panel). N is the 443 number of valid matchups for each parameter. Note that cross validation was not applied for 444 the converted prokaryotes and Prochlorococcus Chl-a because they are the results of the 445 multiplication between their fractions and the TChl-a. As a comparison, the statistics of the 446 previous EOF-based algorithm (without SST) by Xi et al. (2020) for the TChl-a and six PFT 447 Chl-a are also presented (lower panel). Bold marks the improved (or same) statistics.

	Ν	MDPD (%)	RMSD (mg m <sup>-3</sup> )	$\mathbf{R}^2$	MDPDcv (%)	RMSDcv (mg m <sup>-3</sup> )	R <sup>2</sup> cv
EOF-SST hybrid model							
TChl-a	412	30.02	0.85	0.85	30.60	0.88	0.84
Diatoms	296	54.90	0.95	0.79	56.46	1.05	0.78
Dinoflagellates	250	51.44	0.95	0.64	53.81	0.70	0.60
Haptophytes	402	40.24	0.16	0.73	41.96	0.16	0.72
Green algae	285	48.64	0.10	0.59	49.71	0.11	0.56
Prokaryotes	395	44.37	0.09	0.27	44.79	0.09	0.22
f-Prokaryotes	395	40.90	0.20	0.66	41.91	0.20	0.65
Converted Prokaryotes	391	39.36	0.08	0.34	-	-	-
Prochlorococcus	195	41.41	0.02	0.26	44.22	0.02	0.19
f-Prochlorococcus	201	41.08	0.10	0.62	42.07	0.10	0.58
Converted Prochlorococcus	190	39.74	0.02	0.21	-	-	-
Original EOF based algorithm							
TChl-a	394	37.41	1.24	0.76	37.08	1.27	0.75
Diatoms	306	73.70	1.21	0.65	74.74	1.29	0.63
Dinoflagellates	272	55.32	0.93	0.62	57.29	0.72	0.59
Haptophytes	387	47.16	0.22	0.64	48.62	0.24	0.61
Green algae	262	55.81	0.11	0.51	56.26	0.11	0.48
Prokaryotes	367	53.70	0.13	0.15	55.08	0.13	0.11
Prochlorococcus	142	39.65	0.02	0.24	42.68	0.02	0.18





450 Figure 6. Regressions between observed (x-axis, obs.) and predicted (y-axis, pred.) PFT
451 quantities using EOF-SST hybrid algorithm: (A) TChl-a, (B) diatoms, (C) dinoflagellates, (D)
452 haptophytes, (E) green algae, (F) prokaryotes, (G) *f*-prokaryotes, (H) converted prokaryotes

- 453 Chl-a from *f*-prokaryotes, (I) *Prochlorococcus*, (J) *f-Prochlorococcus*, and (K) converted 454 *Prochlorococcus* Chl-a from *f-Prochlorococcus*.
- 455 **3.1.2** SST-separated hybrid algorithms for different SST regimes

456 According to Section 2.2.2, SST-separated hybrid algorithms were performed to retrieve the 457 PFT quantities respectively for the two temperature regimes. Tables 4 summarizes the 458 coefficients fitted in the stepwise regression models based on the whole data set, data set with 459  $SST \ge 8$  °C, and data set with SST < 8 °C, respectively. EOF modes chosen for different 460 PFTs vary with different data sets, and that SST, as an additional regression term, may not 461 always been used in the final prediction models. The term had been identified as insignificant 462 within the stepwise minimization method routine performed in the model regression 463 procedure (Xi et al., 2020). Moreover, the weighting (coefficient) fitted on the SST term changed when different data sets were used. For instance,  $a_{SST}$  fitted in the prediction models 464 465 for prokaryotes using the whole data set was 0.065, and was enhanced to 0.177 for the data set with SST < 8 °C, but was not used in the prediction model for the data set with SST  $\geq 8$ 466 467 °C which is also consistent with the SST-PFT relationship (Figure 3). The 10-point running 468 mean trend showed that SST had a distinct positive correlation with the prokaryotes Chl-a but 469 the correlation turned insignificant when SST was higher than 8 °C.

- 470 As shown in Table 5 and Figure 7, the improvement for TChl-a in the SST-separated 471 algorithms is rather small, indicating that responses in the TChl-a concentration to different 472 SST regimes is relatively stable. For predictions of Chl-a of all PFTs, but Prochlorococcus 473 for which the separation of SST does not apply, the SST-separated algorithms perform much better, indicated by improved statistics in terms of  $R^2$ , RMSD and MDPD (Table 5). 474 475 Prokaryotes, both in quantities of Chl-a and fraction, show the most promising output 476 compared to that from the previous settings. With good performances in retrieving the f-477 prokaryotes and TChl-a, again, the prokaryotes Chl-a is also better derived through converting 478 the fraction to concentrations (Figure 7H) compared to the directly retrieved prokaryotes Chla (Figure 7F), with an increase of  $\mathbb{R}^2$  from 0.34 to 0.43 and reduced MDPD from 39.36% to 479 480 36.51%. Even though prokaryotes retrieval is still not as good as the other PFTs such as 481 diatoms and haptophytes, this is already a significant improvement after a series of 482 experiments by including SST in the retrieval model, establishing separated models based on 483 SST regimes, and retrieving firstly the fraction and performing the conversion.
- 484

Table 4. Coefficients of the EOF-SST based regression models for the merged products used to predict TChl-a, PFT Chl-a and two PFT fractions using data sets with  $SST \ge 8$  °C only, SST < 8 °C only, and all SST, respectively. Note that TChl-a prediction models are based on the EOFs derived from original Rrs spectra but the others are based on the standardized Rrs spectra.

		Ν	Intercept	EOF1	EOF2	EOF3	EOF4	EOF5	EOF6	EOF7	EOF8	SST
	SST≥8°C	353	0.570	154.542	128.794	936.120	-668.923	-275.773	1014.797			-0.017
TChl-a	SST<8°C	59	0.400	91.507	328.418	777.466	-1263.425			-1422.531		
	all SST	412	0.518	143.568	155.376	913.763	-766.036	-225.255	859.576			-0.019
Diatoms	SST≥8°C	239	-5.823	-1.086	1.783	1.346	-1.362	1.713	3.978	2.992		
	SST<8°C	57	-2.351	-0.724	0.922	0.387	-1.621	1.176	3.236			-0.192
	all SST	296	-4.322	-0.961	1.487	1.184	-1.530	1.334	3.598	3.250		-0.070
	SST≥8°C	200	-4.166	-0.561	1.055	0.518	-1.129					-0.032
Dinoflagellates	SST<8°C	50	0.499	1.410	0.578	1.219	-1.834		2.113	7.895		0.085
	all SST	250	-4.753	-0.665	1.111	0.566	-1.206	0.921			-5.563	
Haptophytes	SST≥8°C	343	-4.171	-0.914	1.037		-0.591	-1.279				-0.056
	SST<8°C	59	-3.415	-0.564	0.286	0.522	-1.323	-1.783	-1.930	-3.888		0.049
	all SST	402	-4.983	-1.155	1.061	0.202	-0.714	-0.880		1.314	-2.285	-0.035
	SST≥8°C	244	-2.930		0.706	0.285	-0.566	-1.858	-1.115	2.322		-0.022
Green algae	SST<8°C	41	-3.412		0.478	0.427		-1.484				0.128
	all SST	285	-3.166		0.729	0.235	-0.472	-1.812	-1.631	2.339		
	SST≥8°C	343	-2.410		0.156	0.390		-2.138	-1.039	-1.969		
Prokaryotes	SST<8°C	52	-4.052		0.169	0.687	0.952	-1.561	-2.631	8.019		0.177
	all SST	395	-3.396		0.295	0.379	0.322	-2.090	-2.187	-2.072		0.065
	SST≥8°C	342	0.156	0.703	-0.903	-0.396	0.824	-2.081	-3.082	-4.976		0.032
f-prokaryotes	SST<8°C	53	-1.775	0.689	-0.360		2.061	-2.575	-5.499	8.041		0.124
	all SST	395	-0.945	0.658	-0.698	-0.285	1.034	-1.915	-3.621	-4.659		0.089
Prochlorococcus	SST≥8°C	195	-4.088		0.114			-3.043	-1.619	-3.848		0.028
f-Prochlorococcus	SST≥8°C	201	0.156	0.846	1.184	-0.654	0.564	-2.723		5.067		

491 Table 5. Combined statistics of the regression models from the SST-separated hybrid 492 algorithms for matchup data set with  $SST \ge 8$  °C and that with SST < 8°C. Improved 493 parameters are marked as bold, by comparing to those from the hybrid algorithm without 494 separating SST (Table 3).

	Ν	R <sup>2</sup>	RMSE	<b>MDPD</b> (%)
TChl-a	412	0.86	0.84	30.21
Diatoms	296	0.82	0.93	49.89
Dinoflagellates	250	0.66	0.73	50.23
Haptophytes	402	0.76	0.15	39.53
Green algae	285	0.61	0.10	44.37
Prokaryotes	395	0.35	0.08	40.85
f-Prokaryotes	395	0.72	0.17	33.60
Converted prokaryotes	395	0.43	0.08	36.51

495



497

Figure 7. Combined regressions between observed (x-axis, obs.) and predicted (y-axis, pred.)
PFT quantities from two sets of EOF-SST hybrid algorithms based on different SST ranges:
(A) TChl-a; Chl-a of (B) diatoms, (C) dinoflagellates, (D) haptophytes, (E) green algae, and
(F) prokaryotes; (G) *f*-prokaryotes and (H) converted prokaryotes Chl-a from *f*-prokaryotes.
The dotted black line shows the 1:1 line and the solid black line indicates the regression based
on the whole data set.

## 504 **3.2** Global maps of PFT quantities from merged Rrs products

The improved EOF-SST hybrid algorithms were applied to the merged Rrs and SST products to derive the global TChl-a, PFT Chl-a concentrations and the fractions (see Figure 4, part model application). To illustrate the global distribution of the PFTs, Figures 8-9 show as example the annual mean generated from the derived monthly PFT quantities (except *Prochlorococcus*) for the year of 2011 using the EOF-SST hybrid algorithm (established in Section 2.2.1) and the SST-separated algorithms (established in Section 2.2.2), respectively, 511 with the absolute difference between these two products. Since *Prochlorococcus* barely exist 512 at cold temperatures (cf. Section 2.2), global maps of *Prochlorococcus* Chl-a and f-513 *Prochlorococcus* are generated only for regions with  $SST \ge 8$  °C (Figure 10).

514 In general, distribution patterns of the retrieved TChl-a and the four eukaryotic PFTs from the 515 EOF-SST hybrid algorithm are consistent with that from the combination of the SST-516 separated algorithms (Figure 8). However, distinct differences between the retrievals from the 517 two approaches are found in the polar regions (Figure 8K-O). Compared to that from the 518 EOF-SST hybrid algorithm, TChl-a concentrations derived from SST-separated algorithms 519 are elevated in most regions of the Southern Ocean compared to that from EOF-SST hybrid 520 algorithm, while slightly decreased in the Arctic except for the Greenland Sea and Barents 521 Sea (Figure 8A versus 8F). Diatoms Chl-a from the SST-separated approach is enhanced in 522 the high latitudes and waters near the coasts (Figure 8G versus Figure 8B). Haptophytes Chl-a 523 is also found enhanced mostly in the moderately high latitudes (e.g., 50°S to 60°S near polar 524 fronts) and decreased in very high latitudes and marginal seas using the SST-separated 525 algorithms (Figure N). The Chl-a for dinoflagellates and green algae are generally slightly 526 lower in most of the oceans except for the Southern Ocean near the polar fronts where higher 527 concentrations are observed in the SST-separated algorithms (Figures M&O). Compared to 528 the retrievals of Xi et al., (2020) with detailed inter-comparison among different PFT/PSC 529 products, these satellite derived PFTs from the SST-separated algorithms are in better 530 agreement with the equivalent products from other studies (e.g., Hirata et al., 2011; Losa et 531 al., 2017; Brewin et al., 2015).

532 The retrievals of prokaryotes Chl-a, *f*-prokaryotes and the converted prokaryotes Chl-a using 533 SST-separated algorithms present generally lower Chl-a and fraction for prokaryotes globally 534 except in the regions around 40 °S in the southern hemisphere, 40 °N in the north Pacific 535 Ocean and between 45 °N and 65 °N in the north Atlantic Ocean (Figure 9G-I). The 536 converted Chl-a from fraction shows more reasonable global distribution (Figure 9C and F) 537 compared to the direct retrievals (Figure 9A and D), given the improved performance via 538 fraction conversion for prokaryotes in Section 3.1. However, the EOF-SST hybrid algorithm 539 derived *f*-prokaryotes is saturated (up to 1) in most mid- to low latitudes (Figure 9B). This 540 saturation is remarkably reduced by the SST-separated algorithms (Figure 9E) thanks to the 541 better description of the prokaryotes' dependency on the temperature in the algorithm, 542 resulting in that converted Chl-a from SST-separated algorithms (Figure 9F) has the most 543 reliable retrieval quality among all prokaryotes Chl-a retrievals shown in Figure 9. Compared

to the prokaryotes retrieval in Xi *et al.* (2020), where the original EOF-based algorithm overestimates the prokaryotes Chl-a dramatically, the converted prokaryotes Chl-a in this study shows better agreements with previous retrievals (Hirata *et al.*, 2011; Losa *et al.*, 2017) but is lower in the polar regions. Validation with in situ data is necessary when more measurements are available in the future.

Regarding *Prochlorococcus*, though the fraction converted Chl-a showed no distinct improvement in model performance compared to the direct Chl-a retrieval (Section 3.1.1), the global retrieval depicts the overall decrease in the converted Chl-a by *f-Prochlorococcus* (Figure 10C) compared to the direct retrieval (Figure 10A), which considerably agrees to what we expect from in-situ and other satellite retrievals (Hirata *et al.* 2011, Alvain *et al.* 2008). The conversion to Chl-a is therefore also restricted by the high uncertainty at low TChl-a and also the low variability of *Prochlorococcus* Chl-a (Xi *et al.*, 2020).

In summary, we consider that the SST-separated algorithms perform better than the EOF-SST hybrid algorithm and the original EOF-based algorithm by Xi *et al.* (2020) for global PFT retrievals. Especially for the prokaryotic phytoplankton distributions show more contrasting patterns between gyre and non-gyre regions, and TChl-a and diatom Chl-a are higher in the Southern Ocean where often satellite-derived estimates have shown too low values (e.g., Johnson *et al.*, 2013; Soppa *et al.*, 2014).



Figure 8. Satellite derived estimates of annual (2011) mean surface TChl-a, Chl-a of diatoms, dinoflagellates, haptophytes, and green algae. Panels (A-E): EOF-SST hybrid algorithm with non-separated SST; Panels (F-J): Combined estimates from SST-separated hybrid algorithms for SST  $\geq$  8 °C and SST < 8 °C, respectively. The magenta curve indicates the isotherm of 8 °C. Panels (K-O): Absolute difference between the combined estimates from SST-separated algorithms and that from EOF-SST hybrid algorithm, i.e. panel (F-J) minus panel (A-E).



570 Figure 9. Same as in Figure 8 but for mean surface Chl-a of prokaryotes, *f*-prokaryotes and

571 the converted prokaryotes Chl-a.

569



Figure 10. Satellite derived estimates of annual (2011) mean surface (A) Chl-a of *Prochlorococcus*, (B) *f-Prochlorococcus*, and (C) the converted *Prochlorococcus* Chl-a using
the EOF-SST hybrid algorithm.

576 **3.3 PFT uncertainty** 

# 577 **3.3.1** Look-up table (LUT) for uncertainty due to Rrs

578 Following the steps listed in Section. 2.3.2 to build the look-up table (LUT) for quantifying 579 PFT uncertainty propagated from Rrs uncertainty, Figure 11 and Table 6 show the regressions 580 and the statistical results of  $\sigma_{MC-PFT}$  against the originally predicted PFT quantities. Higher R<sup>2</sup> 581 and lower RMSD are achieved in general when higher degree of polynomial is used (Table 6).

- However, the difference between different regressions is rather small except for diatoms and dinoflagellates, which are caused by their few data points at higher concentrations with corresponding lower MC derived uncertainties (Figures 11B-C). To be conservative, linear regressions were taken as the final LUT functions to determine  $\sigma_{v(Rrs)}$ .
- 586 It is also noted that not for all predictions  $\sigma_{v(Rrs)}$  can be well defined by fitting a (linear) function. For example, the uncertainty from MC simulation for TChl-a varies very little 587 (0.515 - 0.54) (Figure 11A), indicating that the uncertainty is not dependent on the TChl-a 588 concentrations and does not change very much with TChl-a conditions ( $R^2 = 0$ ). Uncertainties 589 of prokaryotic phytoplankton (prokaryotes and Prochlorococcus) Chl-a are not well 590 correlated to their retrievals either ( $R^2 < 0.3$ ); here the uncertainties of their fractions are 591 highly dependent on the fraction retrievals with an inverse correlation ( $R^2 > 0.77$ ); thus the 592 593 retrieval of higher fractions bears lower uncertainty. The derived regressions imply that the LUTs quantify well the PFT uncertainty propagated from Rrs uncertainty for the non-594 595 prokaryotic phytoplankton PFTs Chl-a and the fractions of the prokaryotic phytoplankton but 596 not for the prokaryotic Chl-a. TChl-a uncertainty is relatively stable and not related to the 597 retrieved TChl-a.



Figure 11. Scatterplots of  $\sigma_{PFT}$  based on the MC simulations versus originally retrieved (natural-logarithmic based) PFTs. Regression lines of linear (red), polynomials with degrees of 2 (green) and 3 (blue) fitting are also shown.

Table 6.  $R^2$  and RMSD of the regression functions fitting the relationship between  $\sigma_{PFT}$  and retrieved PFT quantities (using EOF-SST hybrid algorithm).

	L	linear	Poly	nomial 2	Polynomial 3	
	$\mathbb{R}^2$	RMSD	$\mathbf{R}^2$	RMSD	$\mathbb{R}^2$	RMSD
Diatom Chl-a	0.72	0.243	0.73	0.236	0.74	0.235
Dinoflagellates Chl-a	0.71	0.240	0.77	0.213	0.77	0.212
Green algae Chl-a	0.78	0.178	0.78	0.177	0.79	0.177
Haptophytes Chl-a	0.66	0.192	0.66	0.192	0.66	0.191
Prochlorococcus Chl-a	0.28	0.515	0.37	0.482	0.37	0.481
Prokaryotes Chl-a	0.08	0.339	0.10	0.334	0.11	0.332
TChl-a	0.00	0.004	0.01	0.004	0.03	0.004
f-Prochlorococcus	0.84	0.146	0.84	0.146	0.84	0.146
f-prokaryotes	0.77	0.196	0.77	0.196	0.78	0.195

#### 605 **3.3.2 Validation of PFT uncertainty**

The linear LUTs were applied to the retrieved PFT quantities in the testing data set to 606 607 determine the corresponding  $\sigma_{v(Rrs)}$ . The final consolidated uncertainty of the retrieved PFTs 608 from the testing data set,  $\sigma_v$ , were then estimated using Eq. (5). With matchup data of the 609 testing data set (Step 1 of Section 2.3.2), it is possible to assess whether or not the estimated 610 uncertainties for the PFT products are accurate by comparing them to the actual error,  $\delta_{y}$ , defined as  $\delta_y = \ln(c_p) - \ln(c_o)$ . If the uncertainty  $\sigma_y$  is truly representative of its standard 611 612 deviation and thus is reliable, the distribution of the actual errors normalized by the estimated 613 errors  $\delta_v/\sigma_v$  should, to some extent, follow a standard centered normal distribution 614 (Maritorena *et al.*, 2010). Therefore, to validate the estimated PFT uncertainty  $\sigma_v$ , the testing data set was compared to the corresponding  $\delta_y.$  Note that  $\sigma_y$  and  $\delta_y$  are both natural-615 616 logarithmic based.

617 Figure 12 shows the histograms of  $\delta_v / \sigma_v$  distribution derived from the testing data set for all 618 PFT quantities. The corresponding normal distributions determined by the mean and SD are 619 also displayed in comparison with the centered standard normal distribution. For the majority 620 of PFT quantities, the  $\delta_y/\sigma_y$  distribution coincides well with the standard normal distribution, 621 with mean values close to zero and the SD varying from 0.82 to 1.12. Relatively lower SD ( $\leq$ 622 0.70) are found for dinoflagellates Chl-a, Prochlorococcus Chl-a and fraction, and the fraction converted Chl-a both for prokaryotes and Prochlorococcus. Nevertheless, the fraction 623 converted prokaryotes Chl-a presents higher modeled uncertainty (i.e., lower SD of  $\delta_{PFT}/\sigma_{PFT}$ 624 625 in Figure 12H) compared to the direct retrieval of prokaryotes Chl-a (Figure 12F), even though the former has better prediction performance (Table 3). This suggests possible 626 627 underestimation of the actual errors in the direct retrieval for prokaryotes. Prochlorococcus 628 Chl-a and fraction show overall higher modeled uncertainties compared to the actual error 629 (SD  $\leq$  0.67, Figures 12I-K). However, a skewed distribution of  $\delta_y/\sigma_y$  is found for the direct 630 retrieval of *Prochlorococcus* Chl-a (Figure 12I) whereas the *f-Prochlorococcus* and converted 631 Chl-a show little skewness in their  $\delta_v/\sigma_v$ . This suggests the modeled uncertainty is better 632 described for the fraction and the converted Prochlorococcus Chl-a than the direct retrievals. 633 The validation of the PFT uncertainty indicates that our modeled uncertainty is in general 634 close to or higher than the actual error (SD of  $\delta_v/\sigma_v$  is either close to or lower than 1), 635 implying that the uncertainty assessment we performed in this study is reliable and 636 conservative.



638 Figure 12. Distributions of the actual error normalized by the modeled propagated error 639  $(\delta_{PFT}/\sigma_{PFT})$  for all the retrieved (using EOF-SST hybrid algorithm) PFT quantities from the testing data set: (A) TChl-a; Chl-a of (B) diatoms, (C) dinoflagellates, (D) haptophytes (E) 640 green algae, and (F) prokaryotes; (G) *f*-prokaryotes, (H) fraction-converted prokaryotes Chl-a, 641 642 Prochlorococcus Chl-a, (J) (K) **(I)** f-Prochlorococcus, and fraction-converted 643 Prochlorococcus Chl-a. Red and green curves indicate the fitted normal distribution and the 644 standard centered normal distribution, respectively. The red asterisk presents the mean point 645 of the fitted distribution, and the green circle highlights the point of zero. Mean value and 646 standard deviation (SD) of  $\delta_v / \sigma_v$  are also shown.

#### 647 3.3.3 Per-pixel PFT uncertainty

648 Satellite PFT uncertainties were generated for each pixel by applying all uncertainty terms in 649 Eq. (7) to the monthly retrieved PFT products. The annual mean uncertainty for the year of 650 2011 was determined to be consistent with the global PFT maps in Figures 8-10. As a composite product, the annual mean uncertainty was obtained by computing the root mean 651 652 square of the monthly uncertainty product, the same as for the composite uncertainty of other 653 ocean color products, e.g., the ESA Ocean-Colour Climate Change Initiative (OC-CCI) TChl-654 a product (Sathyendranath et al., 2019). It should be noted that in ideal case we should apply 655 the uncertainty quantification scheme firstly to the PFT daily products, and then compute the 656 monthly or yearly composite based on the multi-day uncertainties. However, detailed 657 computations have shown that the monthly composites from the daily products have poorer 658 spatial coverage compared to the directly derived monthly products. This is mainly due to less 659 valid pixels for the nine-band Rrs spectra in the daily to monthly composites than that from 660 the direct monthly products. In addition, the uncertainty derived directly from the monthly products is very comparable to the monthly composite generated based on daily uncertainty 661 (e.g., for TChl-a uncertainty the  $R^2 = 0.98$  and slope = 1.00 between the two derivations, 662 663 details not shown). Therefore, in the present study we applied the uncertainty quantification 664 scheme directly to the monthly PFT products to 1) save computing time and 2) have better 665 spatial coverage for the derived per-pixel uncertainties.

666 Figure 13 demonstrates the annual composite of uncertainties for all the PFT quantities, including also the uncertainty for TChl-a derived from our EOF-SST hybrid algorithm in 667 668 comparison to the uncertainty of the OC-CCI TChl-a product (Figure 13A versus 13D). 669 Uncertainties for diatoms (Figure 13B), dinoflagellates (Figure 13C), haptophytes (Figure 670 13E), and green algae (Figure 13F) display similar distribution pattern with their Chl-a 671 concentration retrievals, with low uncertainties in the gyres, and high in the high latitudes and 672 marginal seas. The overall lowest uncertainty is obtained for haptophytes Chl-a (the natural 673 logarithmic uncertainty varied from 0.23 to 1.26). It is also low for green algae (0.40 - 1.33), 674 whereas diatoms Chl-a show higher uncertainty (0.53-1.73) and dinoflagellates overall the 675 highest (0.78 - 1.70). Regarding the two prokaryotic phytoplankton, the uncertainty for the 676 direct retrieval of prokaryotes Chl-a shows lower uncertainty in the polar regions but higher in 677 the low latitudes (Figure 13G); that for *Prochlorococcus* shows in general high uncertainty in 678 latitudes higher than 40° both south and north and also the gyres (Figure 13J). The 679 uncertainties for the fractions (Figures 13H and 13K) show reverse distributions to the

fraction retrievals with low uncertainty in mid- to low latitudes compared to that in the high latitudes; the uncertainty for *f*-prokaryotes is lower than *f*-*Prochlorococcus* uncertainty. Fraction converted prokaryotes and *Prochlorococcus* Chl-a uncertainties basically follow the patterns of their corresponding fraction uncertainties but are higher (Figures 13I and 13L), because they are derived by combining both the fraction and TChl-a uncertainties, while the latter show little spatial variations (Figure 13A).

686 As the PFT prediction models are based on the multiple linear regressions, the uncertainty of 687 the model coefficients and SST are propagated linearly to the PFT retrievals using Eq. (8). 688 Their corresponding uncertainties are found to have much less spatial variability compared to 689  $\sigma_{v(Rrs)}$ , resulting in that the distribution patterns of the pixel-wise uncertainties generated for 690 the PFT quantities are very much subject to  $\sigma_{v(Rrs)}$  derived from the linear LUT functions. It 691 should be noted that the LUTs built for prokaryotes and Prochlorococcus Chl-a can not 692 represent their uncertainties sufficiently (Figures 11F and 11G), hence their uncertainty 693 products (Figures 13G and 13J) should be used with caution. The uncertainties for the 694 fractions and the converted prokaryotes and Prochlorococcus Chl-a, however, are reported 695 with higher confidence as they are well described by the LUTs and the error propagation 696 analysis. Uncertainties for the other PFT quantities are in general well and conservatively 697 quantified, which are also justified by the uncertainty validation. Yet uncertainties for satellite 698 derived PFTs have been rarely reported, except that Brewin, Ciavatta, et al. (2017) performed 699 uncertainty evaluation on phytoplankton size classes and two phytoplankton groups (diatoms 700 and dinoflagellates) retrieved by the re-tuned abundance/ecological based algorithm in the 701 North Atlantic region. Though the exact values of the uncertainty estimates are not provided 702 by Brewin, Ciavatta, et al. (2017), it is seen that their uncertainty maps for diatoms and 703 dinoflagellates are within the same order of magnitudes (when converted from the base-10 704 logarithm to the natural-logarithm) and show similar distribution patterns with our uncertainty 705 products in the North Atlantic Ocean.

Compared to the reported uncertainty for the TChla OC-CCI product, the uncertainty of TChl-a from our EOF-SST hybrid algorithm shows much lower uncertainty. Our TChl-a uncertainty varies from 0.53 to 0.58 (Figure 13A), presenting very little spatial variability with only slightly higher uncertainty in the gyres and some marginal seas. The OCCCI chlorophyll uncertainty ranged between 0.43 and 1.18 (Figure 13D), showing large spatial variability with much higher uncertainty in the marginal seas and high latitudes compared to our TChl-a uncertainty. Low uncertainty in and surrounding the gyres is comparable with our uncertainty estimates. Remarkably reduced uncertainty of the TChl-a derived by our algorithm in high latitudes (> 40°) indicates that the EOF-SST hybrid algorithm has great potential in improving TChl-a estimation especially in polar regions where the standard OC algorithms always introduce high errors (IOCCG, 2015). The overall lower uncertainty in our TChl-a product also reveals at a certain scale that our estimates for most of the PFT quantities and the corresponding uncertainties are reliable.



719

Figure 13. Per-pixel uncertainty (in natural logarithmic scale) of the annual mean of 2011 for the satellite derived PFT quantities from EOF-SST hybrid algorithm. The dashed-line box frames in particular (A) the uncertainty of TChl-a from the EOF-SST hybrid algorithm in comparison to (D) the OC-CCI TChl-a uncertainty.

### 724 **4 Summary**

This study improved the previously established EOF-based approach by Xi *et al.* (2020) for estimating globally the Chl-a of six PFTs using merged ocean color Rrs products. The modified retrieval scheme, named EOF-SST hybrid algorithm, was developed by using updated input data sets and accounting for the influence of SST on different PFT quantities. 729 Furthermore, fractions of prokaryotes and *Prochlorococcus* were also included as retrieved 730 PFT quantities, which lead to more accurate retrievals compared to their Chl-a retrievals. The 731 fraction retrievals were also used together with TChl-a retrievals to obtain a fraction-732 converted Chl-a for the two PFTs. The latter shows prominent improvements for prokaryotes 733 Chl-a, but not for *Prochlorococcus* Chl-a. By further splitting the input data set according to 734 the PFT dependence on SST in different SST regimes, separated retrieval algorithms for low 735 and high temperature waters were established, presenting even much more improved 736 performance for all PFTs than the hybrid algorithm based on the whole data set. 737 Improvements for PFT retrievals were mostly obtained in the high latitudes. Finally, the 738 pixel-by-pixel uncertainty of the satellite PFT retrievals was assessed by accounting for the 739 uncertainties from input data and model parameters via an error propagation method. These 740 satellite PFT uncertainties, for the first time reported on global scale and for spectral-based 741 PFT retrieval approaches, provide reliable error estimates for the PFT products which allow 742 us to better understand the product quality both in time and space.

743 This study uses the GlobColour merged OC products that span a period from 2002 to 2012 744 only. However, our EOF-SST hybrid algorithm including pixelwise uncertainties can easily 745 be expanded to other OC sensors, such as MODIS-VIIRS merged and OLCI products. 746 Uncertainty assessment for the PFT estimates from different satellite products, is needed for 747 consistent long-term PFT data set from multiple satellites with assured continuity. Such a data 748 set is required to enable tracking the shifting in phytoplankton community structure under the 749 changing climate for example. PFT products with uncertainty estimates are also beneficial to 750 applications of ecosystem modelling by helping to simulate and/or evaluate the model outputs, 751 as well as being assimilated within these models to further improve forecasting marine 752 biogeochemistry, as used by many marine services.

### 753 Acknowledgement

754 This work is supported by a collaborative project, OLCI-PFT (ACRI-AWI Offer #209-755 180104), between ACRI-ST and Phytooptics team at Alfred Wegener Institute. The 756 contribution by S.L. was partially funded by the Deutsche Forschungsgemeinschaft (DFG, 757 German Research Foundation)-Project number 268020496-TRR 172, within the 758 Transregional Collaborative Research Center ArctiC Amplication: Climate Relevant 759 Atmospheric and SurfaCe Processes, and Feedback Mechanisms (AC)3 (Project C03), and 760 was also partly made in the framework of the state assignment of the Federal Agency for 761 Scientific Organizations (FASO) Russia (theme 0149-2019-0015). We thank Marc Taylor for 762 the primary R script on the EOF method, Sonja Wiegmann for her dedication in measuring 763 and processing the HPLC data for our own cruises, and all the previous and current 764 Phytooptics team members who participated in the past cruises for data collection and 765 analysis. We are thankful to all the scientists and crew involved in the global HPLC data 766 collection and analyses for providing their pigment data. Thanks also to NASA, ESA and 767 EUMETSAT for the SeaWiFS, MODIS, and MERIS data, and specially the GlobColour program for providing the merged ocean color L3 products. GlobColour data 768 769 (http://globcolour.info) used in this study has been developed, validated, and distributed by ACRI-ST, France. This study has been conducted using E.U. Copernicus Marine Service 770 771 Information in terms of CMEMS sea surface temperature products. We also thank the Ocean-772 Colour Climate Change Initiative (OC-CCI) program for providing global chlorophyll-a data 773 products.

### 774 Data availability

The DPA derived PFT Chl-a for diatoms, haptophytes and prokaryotes from the pigment database I were published already in Losa *et al.* (2017) and are available from PANGAEA: https://doi.pangaea.de/10.1594/PANGAEA.875879 (Soppa *et al.*, 2017). All matchup data used in this study, including the collocated *in situ* data, nine-band Rrs data from the GlobColour merged products, and SST data, are also to be prepared and submitted to PANGAEA by end of January 2021.

### 781 **References**

- ACRI-ST GlobColour Team, Mangin, A., & Hembise Fanton d'Andon, O. (2017). *GlobColour Product User Guide* (GC-UM-ACR-PUG-01, Version 4.1). Sophia-Antipolis: ACRI-ST.
- 784 Aiken, J., Pradhan, Y., Barlowd, R., Lavender, S., Poulton, A., Patrick, H., & Hardman-Mountford, N. 785 (2009). Phytoplankton pigments and functional types in the Atlantic Ocean: A decadal 786 assessment, 1995-2005. Deep Sea Res., Part 56(15), 899-917. II. 787 https://doi.org/10.1016/j.dsr2.2008.09.017
- Alvain, S., Moulin, C., Dandonneau, Y., & Bréon, F. M. (2005). Remote sensing of phytoplankton
  groups in case 1 waters from global SeaWiFS imagery. *Deep Sea Res.*, *Part I*, 52(11), 1989–
  2004. https://doi.org/10.1016/j.dsr.2005.06.015
- Alvain, S., Moulin, C., Dandonneau, Y., & Loisel, H. (2008). Seasonal distribution and succession of
   dominant phytoplankton groups in the global ocean: A satellite view. *Global Biogeochem. Cycles*, 22(3), GB3001, https://doi.org/10.1029/2007GB003154
- Amante, C., & Eakins, B.W. (2009). *ETOPO1 1 Arc-Minute Global Relief Model: Procedures, Data Sources and Analysis* (NOAA Technical Memorandum NESDIS NGDC-24). Boulder, CO:
   National Geophysical Data Center.
- Antoine, D., d'Ortenzio, F., Hooker, S. B., Bécu, G., Gentili, B., Tailliez, D., & Scott, A. J. (2008).
  Assessment of uncertainty in the ocean reflectance determined by three satellite ocean color sensors (MERIS, SeaWiFS and MODIS-A) at an offshore site in the Mediterranean Sea

- 800
   (BOUSSOLE project).
   J
   Geophys.
   Res.
   Ocean.,
   113,
   C07013.

   801
   https://doi.org/10.1029/2007JC004472

   </td
- 802 Bracher, A., Xi, H., Dinter, T., Mangin, A., Strass, V. H., von Appen, W.-J., & Wiegmann, S. (2020). 803 High resolution water column phytoplankton composition across the Atlantic Ocean from ship-804 7, undulating radiometry. Front. towed vertical Mar. Sci., 235. 805 https://doi.org/10.3389/fmars.2020.00235
- Bracher, A., Bouman, H. A., Brewin, R. J. W., Bricaud, A., Brotas, V., Ciotti, A. M., et al. (2017).
  Obtaining Phytoplankton Diversity from Ocean Color: A Scientific Roadmap for Future
  Development. *Front. Mar. Sci.*, 4, 1–15. https://doi.org/10.3389/fmars.2017.00055
- Bracher, A., Taylor, M.H., Taylor, B., Dinter, T., Röttgers, R., & Steinmetz, F. (2015). Using
  empirical orthogonal functions derived from remote-sensing reflectance for the prediction of
  phytoplankton pigment concentrations. *Ocean Sci.*, 11, 139-158. https://doi.org/10.5194/os-11139-2015
- Bracher, A., Vountas, M., Dinter, T., Burrows, J. P., Röttgers, R., & Peeken, I. (2009). Quantitative
  observation of cyanobacteria and diatoms from space using PhytoDOAS on SCIAMACHY data. *Biogeosciences*, 6, 751-764. https://doi.org/10.5194/bg-6-751-2009
- Brewin, R. J. W., Ciavatta, S., Sathyendranath, S., Jackson, T., Tilstone, G., Curran, K., et al. (2017).
  Uncertainty in Ocean-Color Estimates of Chlorophyll for Phytoplankton Groups. *Front. Mar. Sci.*, *4*, 104. https://doi.org/10.3389/fmars.2017.00104
- Brewin, R. J. W., Tilstone, G., Jackson, T., Cain, T., Miller, P., Lange, P. K., et al. (2017). Modelling
  size-fractionated primary production in the Atlantic Ocean from remote sensing. *Prog. Oceanogr.*, 158, 130-149. https://doi.org/10.1016/j.pocean.2017.02.002
- Brewin, R. J. W., Sathyendranath, S., Hirata, T., Lavender, S. J., Barciela, R. M., & HardmanMountford, N. J. (2010). A three-component model of phytoplankton size class for the Atlantic
  Ocean. *Ecol. Modell.*, 221(11), 1472-1483. https://doi.org/10.1016/j.ecolmodel.2010.02.014
- Brewin, R. J. W., Sathyendranath, S., Jackson, T., Barlow, R., Brotas, V., Airs, R., & Lamont, T.
  (2015). Influence of light in the mixed-layer on the parameters of a three-component model of
  phytoplankton size class. *Remote Sens. Environ.*, 168, 437-450.
  https://doi.org/10.1016/j.rse.2015.07.004
- Brotas, V., Brewin, R. J. W., Sá, C., Brito, A. C., Silva, A., Mendes, C. R., et al. (2013). Deriving
  phytoplankton size classes fromsatellite data: validation along a trophic gradient in the eastern
  Atlantic Ocean. *Remote Sens. Environ.*, 134, 66–77. https://doi.org/10.1016/j.rse.2013.02.013
- Ciotti, A. M., & Bricaud, A. (2006). Retrievals of a size parameter for phytoplankton and spectral light
  absorption by colored detrital matter from water-leaving radiances at SeaWiFS channels in a
  continental shelf region off Brazil. *Limnol. Oceanogr. Methods*, 4, 237–253.
  https://doi.org/10.4319/lom.2006.4.237
- 836 Claustre, H., Hooker, S. B., Van Heukelem, L., Berthon, J.-F., Barlow, R., Ras, J., et al. (2004). An 837 intercomparison of HPLC phytoplankton pigment methods using in situ samples: Application to 838 sensing and database activities. Mar. Chem., 85(1-2), 41-61. remote 839 https://doi.org/10.1016/j.marchem.2003.09.002
- Correa-Ramirez, M., Morales, C. E., Letelier, R., Anabalon, V., & Hormazabal, S. (2018). Improving
  the remote sensing of phytoplankton functional types (PFT) using empirical orthogonal
  functions: a case study in a coastal upwelling region. *Remote Sens.*, 10(4), 498.
  https://doi.org/10.3390/rs10040498
- Craig, S. E., Jones, C. T., Li, W. K. W., Lazin, G., Horne, E., Caverhill, C., & Cullen, J. J. (2012).
  Deriving optical metrics of coastal phytoplankton biomass from ocean colour. *Remote Sens. Environ.*, 119, 72–83. https://doi.org/10.1016/j.rse.2011.12.007

- de Mora, L., Butenschön, M., & Allen, J. I. (2016). The assessment of a global marine ecosystem
  model on the basis of emergent properties and ecosystem function: a case study with ERSEM. *Geosci. Model Dev.*, 9, 59–76. https://doi.org/10.5194/gmd-9-59-2016
- Bevred, E., Sathyendranath, S., & Platt, T. (2009). Decadal changes in ecological provinces of the
  Northwest Atlantic Ocean revealed by satellite observation. *Geophys. Res. Letter.*, *36*, L19607.
  https://doi.org/10.1029/2009GL039896
- Bevred, E., Sathyendranath, S., Stuart, V., Maass, H., Ulloa, O., & Platt, T. (2006). A two-component
  model of phytoplankton absorption in the open ocean: Theory and applications. J. *Geophys. Res. Ocean.*, 111, C03011, https://doi.org/10.1029/2005JC002880
- Bonlon, C. J., Martin, M., Stark, J., Roberts-Jones, J., Fiedler, E., & Wimmer, W. (2012). The
  operational sea surface temperature and sea ice analysis (OSTIA) system. *Remote Sens. Environ.*, 116, 140-158. https://doi.org/10.1016/j.rse.2010.10.017
- Falkowski, P. G., Laws, E. A., Barber, R. T., & Murray, J. W. (2003). Phytoplankton and their role in
  primary, new, and export production. In M.J.R. Fasham (Eds.), *Ocean biogeochemistry. Global Change The IGBP Series (closed)* (pp99-121). Berlin, Heidelberg: Springer.
  https://doi.org/10.1007/978-3-642-55844-3\_5
- Flombaum, P., Gallegos, J. L., Gordillo, R. A., Rincon, J., Zabala, L. L., Jiao, N., et al. (2013). Present
  and future global distributions of the marine Cyanobacteria Prochlorococcus and Synechococcus. *Proc. Natl. Acad. Sci. 110*, 9824–9829. https://doi.org/10.1073/pnas.1307701110
- Hirata, T., Hardman-Mountford, N. J., Brewin, R. J. W., Aiken, J., Barlow, R., Suzuki, K., et al.
  (2011). Synoptic relationships between surface Chlorophyll-a and diagnostic pigments specific
  to phytoplankton functional types. *Biogeosciences*, 8, 311–327. https://doi.org/10.5194/bg-8-311-2011
- 870 Hirata, T., Saux Picart, S., Hashioka, T., Aita-Noguchi, M., Sumata, H., Shigemitsu, M., et al. (2013). 871 A comparison between phytoplankton community structures derived from a global 3D ecosystem 872 model satellite and observation. J. Mar. Syst., 109–101, 129–137. 873 https://doi.org/10.1016/j.jmarsys.2012.01.009
- 874 Holt, J., Allen, J. I., Anderson, T. R., Brewin, R. J. W., Butenschön, M., Harle, J., et al. (2014). 875 Challenges in integrative approaches to modelling the marine ecosystems of the North Atlantic: 876 fish and coasts to ocean. Oceanogr., 129. physics to Prog. 285–313. 877 https://doi.org/10.1016/j.pocean.2014.04.024
- 878 IOCCG (2014). Phytoplankton functional types from space. In S. Sathyendranath and V. Stuart (Eds.)
   879 *Reports of the International Ocean Color Coordinating Group* (No. 15). Dartmouth, NS:
   880 IOCCG.
- IOCCG (2015). Ocean Colour Remote Sensing in Polar Seas. In M. Babin, K. Arrigo, S. Bélanger,
   M.-H. Forget (Eds.) *Reports of the International Ocean Color Coordinating Group* (No. 16).
   Dartmouth, NS: IOCCG.
- IOCCG (2019). Uncertainties in ocean colour remote sensing. In F. Melin (Eds.) *Reports of the International Ocean Color Coordinating Group* (No. 18). Dartmouth, NS: IOCCG.
- JCGM (2008). Evaluation of Measurement Data Guide to the Expression of Uncertainty in
   Measurement. JCGM 100:2008.
- Kostadinov, T. S., Milutinović, S., Marinov, I., & Cabré, A. (2016). Carbon-based phytoplankton size
  classes retrieved via ocean color estimates of the particle size distribution. *Ocean Sci.*, *12*, 561–
  575. https://doi.org/10.5194/os-12-561-2016
- Kostadinov, T. S., Siegel, D. A., & Maritorena, S. (2009). Retrieval of the particle size distribution
  from satellite ocean color observations. J. Geophys. Res. Ocean., 114, 09015.
  https://doi.org/10.1029/2009JC005303

- Lange P. K., Werdell P. J., Erickson Z. K., Dall'Olmo G., Brewin R., Zubkov M., et al. (2020).
  Radiometric approach for the detection of picophytoplankton assemblages across oceanic fronts. *Optics Express.*, 28(18), 25682-25705. https://doi.org/10.1364/OE.398127
- Le Quéré, C., Harrison, S. P., Prentice, C. I., Buitenhuis, E. T., Aumonts, O., Bopp L., et al. (2005).
  Ecosystem dynamics based on plankton functional types for global ocean biogeochemistry
  models. *Global Change Biology*, *11*, 2016-2040. https://doi.org/10.1111/j.13652486.2005.1004.x
- Lee, Z., Arnone, R., Hu, C., Werdell, P. J., & Lubac, B. (2010). Uncertainties of optical parameters
  and their propagations in an analytical ocean color inversion algorithm. *Appl. Optics*, 49, 369–
  381. https://doi.org/10.1364/AO.49.000369
- Losa, S. N., Soppa, M. A., Dinter, T., Wolanin, A., Brewin, R.J.W., Bricaud, A., et al. (2017).
  Synergistic Exploitation of Hyper- and Multi-Spectral Precursor Sentinel Measurements to
  Determine Phytoplankton Functional Types (SynSenPFT). *Front. Mar. Sci.*, 4, 1–22.
  https://doi.org/10.3389/fmars.2017.00203
- Losa, S. N., Dutkiewicz, S., Losch, M., Oelker, J., Soppa, M. A., Trimborn, S., Xi, H., & Bracher, A.
  (2019). On modeling the Southern Ocean Phytoplankton Functional Types, *Biogeosciences Discuss.*, https://doi.org/10.5194/bg-2019-289
- Lubac, B., & Loisel, H. (2007). Variability and classification of remote sensing reflectance spectra in
  the eastern English Channel and southern North Sea. *Remote Sens. Environ.*, *110*(1), 45-58.
  https://doi.org/10.1016/j.rse.2007.02.012
- Maritorena, S., Hembise Fanton d'Andon, O., Mangin, A., & Siegel, D. A. (2010). Merged satellite
  ocean color data products using a bio-optical model: characteristics, benefits and issues. *Remote Sens. Environ.* 114, 1791–1804. https://doi.org/10.1016/j.rse.2010.04.002
- McKinna, L. I. W., Cetinic, I., Chase, A. P. & Werdell, P. J. (2019). Approach for Propagating
  Radiometric Data Uncertainties Through NASA Ocean Color Algorithms. *Front. Earth Sci.*, 7,
  176. https://doi.org/10.3389/feart.2019.00176
- Mélin, F., & Franz, B. A. (2014). Assessment of satellite ocean colour radiometry and derived
  geophysical products. In G. Zibordi, C. Donlon, A. Parr (Eds.) *Optical Radiometry for Oceans Climate Measurements*, Experimental Methods in the Physical Sciences (Vol. 47). Academic
  Press.
- Mouw, C., Hardman-Montford, N., Alvain, S., Bracher, A., Brewin, R. J. W., Bricaud, A., et al.
  (2017). A Consumer's Guide to Satellite Remote Sensing of Multiple Phytoplankton Groups in the Global Ocean. *Front. Mar. Sci.*, *4*, 00041, https://doi.org/10.3389/fmars.2017.00041
- Johnson, R., Strutton, P. G., Wright, S. W., McMinn, A., & Meiners, K. M. (2013). Three improved
  satellite chlorophyll algorithms for the Southern Ocean. J. Geophy. Res. Ocean., 118, 3694-3703.
  https://doi.org/10.1002/jgrc.20270
- Palacz, A. P., St. John, M. A., Brewin, R. J. W., Hirata, T., & Gregg, W. W. (2013). Distribution of
  phytoplankton functional types in high-nitrate, low-chlorophyll waters in a new diagnostic
  ecological indicator model. *Biogeosciences*, 10, 7553-7574. https://doi.org/10.5194/bg-10-75532013
- Pradhan, H. K., Völker, C., Losa, S. N., Bracher, A., & Nerger L. (2019). Assimilation of global total
  chlorophyll OC-CCI data and its impact on individual phytoplankton fields. J. Geophys. Res.
  Ocean., 124, 470-490. https://doi.org/10.1029/2018JC014329
- Pradhan, H. K., Völker, C., Losa, S. N., Bracher, A., & Nerger, L. (2020). Global assimilation of
  ocean-color data of phytoplankton functional types: Impact of different datasets. *J. Geophys. Res. Ocean.*, 125, e2019JC015586. https://doi.org/10.1029/2019JC015586

- 940 Qi, L., Lee, Z., Hu, C., & Wang, M. (2017). Requirement of minimal signal-to-noise ratios of ocean
  941 color sensors and uncertainties of ocean color products. *J. Geophys. Res. Ocean.*, 122, 2595–
  942 2611. https://doi.org/10.1002/2016JC012558
- Raitsos, D. E., Lavender, S. J., Maravelias, C. D., Haralabous, J., Richardson, A. J., & Reid, P. C.
  (2008). Identifying four phytoplankton functional types from space: An ecological approach. *Limnol. Oceanogr.*, 53(2), 605-613. https://doi.org/10.4319/lo.2008.53.2.0605
- Sathyendranath, S., Brewin, R. J. W., Brockmann, C., Brotas, V., Calton, B., Chuprin, A., et al. (2019).
  An ocean-colour time series for use in climate studies: the experience of the Ocean-Colour Climate Change Initiative (OC-CCI). *Sensors*, *19*, 4285. https://doi.org/10.3390/s19194285
- 949Sathyendranath, S., Jackson, T., Brockmann, C., Brotas, V., Calton, B., Chuprin, A., et al. (2020).950ESA Ocean Colour Climate Change Initiative (Ocean\_Colour\_cci): Global chlorophyll-a data951products gridded on a sinusoidal projection, Version 4.2. Centre for Environmental Data952Analysis, 210ctober2020
- 953 *citation*. https://catalogue.ceda.ac.uk/uuid/99348189bd33459cbd597a58c30d8d10
- Soja-Woźniak, M., Craig, S. E., Kratzer, S., Wojtasiewicz, B., Darecki, M., & Jones, C. T. (2017). A
  Novel Statistical Approach for Ocean Colour Estimation of Inherent Optical Properties and
  Cyanobacteria Abundance in Optically Complex Waters. *Remote Sens.*, 9, 1–22.
  https://doi.org/10.3390/rs9040343
- Soppa, M. A., Hirata, T., Silva, B., Dinter, T., Peeken, I., Wiegmann, S., & Bracher, A. (2014). Global
  retrieval of diatom abundance based on phytoplankton pigments and satellite data. *Remote Sens.*6, 10089–10106. https://doi.org/10.3390/rs61010089
- Soppa, M.A., Peeken, I., & Bracher, A. (2017). Global chlorophyll "a" concentrations for diatoms,
  haptophytes and prokaryotes obtained with the diagnostic pigment analysis of HPLC data
  compiled from several databases and individual cruises. PANGAEA.
  https://doi.org/10.1594/PANGAEA.875879
- Taylor, B. B., Taylor, M. H., Dinter, T., & Bracher, A. (2013). Estimation of relative phycoerythrin
   concentrations from hyperspectral underwater radiance measurements A statistical approach. J.
   *Geophys. Res. Ocean.* 118, 2948–2960. https://doi.org/10.1002/jgrc.20201
- 968 Uitz, J., Claustre, H., Morel, A., & Hooker, S. B. (2006). Vertical distribution of phytoplankton
  969 communities in open ocean: An assessment based on surface chlorophyll. J. *Geophys. Res.*970 *Ocean.*, 111, C08005. https://doi.org/10.1029/2005JC003207
- Vidussi, F., Claustre, H., Manca, B. B., Luchetta, A., & Marty, J.-C. (2001). Phytoplankton pigment
  distribution in relation to upper thermocline circulation in the eastern Mediterranean Sea during
  winter. J. Geophys. Res. Ocean., 106(C9), 19,939-19,956. https://doi.org/10.1029/1999JC000308
- Ward, B. A., Dutkiewicz, S., Jahn, O., & Follows, M. J. (2012). A size-structured food-web model for
  the global ocean. *Limnol. Oceanogr.* 57, 1877–1891. https://doi.org/10.4319/lo.2012.57.6.1877
- Ward, B.A. (2015). Temperature-correlated changes in phytoplankton community structure are
   restricted to polar waters. *PLoS One*. https://doi.org/10.1371/journal.pone.0135581
- Werdell, P. J., McKinna, L., Boss, E., Ackleson, S., Craig, S., Gregg, W., et al. (2018). An overview
  of approaches and challenges for retrieving marine inherent optical properties from ocean color
  remote sensing. *Prog. Oceanogr.*, *160*, 186-212. https://doi.org/10.1016/j.pocean.2018.01.001
- Werdell, P. J., Roesler, C. S., & Goes, J. I. (2014). Discrimination of phytoplankton functional groups
  using an ocean reflectance inversion model. *Appl. Opt.*, 53, 4833.
  https://doi.org/10.1364/ao.53.004833
- Worsfold, M., Good, S., McLaren, A., Fiedler, E., Roberts-Jones, J., & Martin, M. (2014). Quality
   Information Document Global Ocean OSTIA Sea Surface Temperature Reprocessing SST-GLO SST-L4-REP-OBSERVATIONS-010-011. Ref: CMEMS-SST-QUID-010-011.

- Xi, H., Losa, S. N., Mangin, A., Soppa, M. A., Garnesson, P., Demaria, J., et al. (2020). Global retrieval of phytoplankton functional types based on empirical orthogonal functions using CMEMS GlobColour merged products and further extension to OLCI data. *Remote Sens. Environ.*, 240, 111704. https://doi.org/10.1016/j.rse.2020.111704
- Xiao, Y., & Friedrichs, M. A. M. (2014). The assimilation of satellite-derived data into a onedimensional lower trophic level marine ecosystemmodel. J. Geophys. Res. Ocean., 119, 2691–
  2712. https://doi.org/10.1093/plankt/fbp127