



**HAL**  
open science

## Assessment of a regional physical–biogeochemical stochastic ocean model. Part 2: Empirical consistency

Vassilios D Vervatis, Pierre de Mey-Frémaux, Nadia Ayoub, John Karagiorgos, Stefano Ciavatta, Robert J W Brewin, Sarantis Sofianos

### ► To cite this version:

Vassilios D Vervatis, Pierre de Mey-Frémaux, Nadia Ayoub, John Karagiorgos, Stefano Ciavatta, et al.. Assessment of a regional physical–biogeochemical stochastic ocean model. Part 2: Empirical consistency. *Ocean Modelling*, 2021, 160, pp.101770. 10.1016/j.ocemod.2021.101770 . hal-03441458

**HAL Id: hal-03441458**

**<https://hal.science/hal-03441458>**

Submitted on 22 Nov 2021

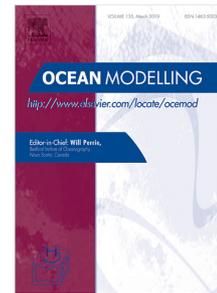
**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## Journal Pre-proof

Assessment of a regional physical-biogeochemical stochastic ocean model. Part 2: Empirical consistency

Vassilios D. Vervatis, Pierre De Mey-Frémaux, Nadia Ayoub,  
John Karagiorgos, Stefano Ciavatta, Robert J.W. Brewin,  
Sarantis Sofianos



PII: S1463-5003(21)00020-2  
DOI: <https://doi.org/10.1016/j.ocemod.2021.101770>  
Reference: OCEMOD 101770

To appear in: *Ocean Modelling*

Received date : 18 May 2020  
Revised date : 19 January 2021  
Accepted date : 7 February 2021

Please cite this article as: V.D. Vervatis, P. De Mey-Frémaux, N. Ayoub et al., Assessment of a regional physical-biogeochemical stochastic ocean model. Part 2: Empirical consistency. *Ocean Modelling* (2021), doi: <https://doi.org/10.1016/j.ocemod.2021.101770>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2021 Published by Elsevier Ltd.

## 1 Assessment of a regional physical-biogeochemical stochastic ocean model. Part 2: 2 empirical consistency

3 Vassilios D. Vervatis (1), Pierre De Mey-Frémaux (2), Nadia Ayoub (2), John Karagiorgos (1),  
4 Stefano Ciavatta (3, 4), Robert J. W. Brewin (3, 5), and Sarantis Sofianos (1)

5 (1) University of Athens, Department of Physics, Athens, Greece. (2) LEGOS/CNRS,  
6 Toulouse, France. (3) Plymouth Marine Laboratory, Plymouth, UK. (4) National Centre for  
7 Earth Observation, Plymouth, UK. (5) University of Exeter, College of Life and Environmental  
8 Sciences, Penryn, Cornwall, UK.

9 *Correspondence to:* Vassilios D. Vervatis (vervatis@oc.phys.uoa.gr)

10 **Abstract.** In this Part 2 article of a two-part series, observations based on satellite missions  
11 were used to evaluate the empirical consistency of model ensembles generated via stochastic  
12 modelling of ocean physics and biogeochemistry. A high-resolution Bay of Biscay  
13 configuration was used as a case study to explore the model error subspace in both the open  
14 and coastal ocean. In Part 1 of this work, three experiments were carried out to generate model  
15 ensembles by perturbing only physics, only biogeochemistry, and both of them simultaneously.  
16 In Part 2 of this work, empirical consistency was checked, first by means of rank histograms  
17 projecting the data onto the model ensemble classes, and second, by pattern-selective  
18 consistency criteria in the space of “array modes” (eigenvectors of the representer matrix).  
19 Rank histograms showed large dependency on geographical region and on season for sea  
20 surface temperature (SST), sea-level anomaly (SLA), and phytoplankton functional types  
21 (PFT), shifting from consistent model-data configurations to large biases because of model  
22 ensemble underspread. Consistency for SST array modes was found to be verified at large,  
23 small and coastal scales soon after the ensemble spin-up. Array modes for the along-track sea-  
24 level showed useful consistent information at large scales and at the mesoscale; for the gridded  
25 SLA was verified only at large scale. Array modes showed that biogeochemical model  
26 uncertainties generated by stochastic physics, were effectively detected by PFT measurements  
27 at large scales, as well as at mesoscale and small-scale. By contrast, perturbing only  
28 biogeochemistry, with an identical physical forcing across the ensemble, limits the potential of  
29 PFT measurements at detecting and possibly correcting small-scale biogeochemical model  
30 errors. When an ensemble was found to be inconsistent with observations along a particular  
31 direction (here, an array mode), a plausible reason is that other error processes must have been  
32 active in the model, in addition to the ones at work across the ensemble.

33 **Keywords:** stochastic modelling, ensembles, phytoplankton functional types, prior error  
34 covariances, array modes, Bay of Biscay

### 35 1 Introduction

36 In the coastal parts of regional ocean models, many factors can complicate the assimilation of  
37 data. One of them is the characterization and specification of model errors, which are critical  
38 in any assimilation scheme, but extremely challenging in the coastal zone. Ocean model errors  
39 strongly depend on spatiotemporal scales, though any attempt at separation is confounded by  
40 strong nonlinearity in the dynamics that can couple variability at different frequencies and  
41 wavenumbers (*Auclair et al.*, 2003). Another factor complicating data assimilation is the  
42 specification of observational errors, which are made up of measurement (usually small) and  
43 representativity errors (usually large or unknown), and cross-correlations between observations  
44 (usually unknown) (*Oke and Sakov*, 2008).

45 Most studies point at the benefit of advanced assimilation methods with built-in error  
46 propagation (*Kourafalou et al.*, 2015a; 2015b), such as the Ensemble Kalman Filter (EnKF;

47 *Evensen*, 2003) and variants. As a first step, one must characterize the forecast uncertainties  
48 under various error regimes (situations where the model is in error), either in response to time-  
49 varying forcing errors (e.g. uncertainties in boundary conditions; *Ghantous et al.*, 2020) or due  
50 to internal sources (e.g. model parameterizations; *Brankart et al.*, 2015), and include realistic  
51 error dynamics through stochastic modelling. For instance, depending on the wind regime, or  
52 seasonal baroclinic instabilities of the slope current, shelf errors can be different in terms of  
53 spatiotemporal scales (*Auclair et al.*, 2003). Biogeochemical model errors can also be different  
54 stemming from unresolved scales and biodiversity (*Garnier et al.*, 2015). It should be noted  
55 that stochastic modelling, in itself, may not yield realistic error dynamics, for example, it could  
56 inflate the ensemble (*Anderson*, 2009).

57 In the companion article Part 1 (*Vervatis et al.*, 2021), we configured a high-resolution ( $1/36^\circ$ )  
58 stochastic ocean model for the Bay of Biscay performing physical-biogeochemical ensemble  
59 simulations. We carried out quantitative assessment of model-data misfits and qualitative  
60 evaluation of multivariate incremental analysis. We found that the skill of the perturbation  
61 method to generate model errors was improved (in general) for physics compared with  
62 biogeochemical perturbations in source and sink terms (i.e. larger ensemble spread when  
63 physics is perturbed) and that the data assimilation performance to correct those model errors,  
64 was largely dependent on the chosen multivariate analysis. On the other hand, the  
65 biogeochemical model spread was found under-dispersive, leading to disjoint supports<sup>1</sup> of the  
66 model and data probability density functions (*pdf*), and performance was mainly defined by  
67 the assimilation of ocean colour data, possibly because of weak cross-covariances between  
68 ocean physics and biogeochemistry.

69 In light of these findings, in this companion article Part 2, we present two methods to evaluate  
70 the consistency with respect to observations (hereafter empirical consistency) of a stochastic  
71 coupled ocean model, which comes in the form of an ensemble of multivariate ocean states.  
72 We focus on observational products including physics and biogeochemistry, derived from  
73 satellite missions monitoring upper-ocean properties, such as the sea surface temperature  
74 (SST), the sea-level anomaly (SLA) and an ocean colour product of phytoplankton functional  
75 types (PFT; *Brewin et al.*, 2010; 2015; 2017). Several biogeochemical studies have focused on  
76 improving the model's predictive skill using remote sensing and in-situ observations,  
77 incorporating data assimilation and probabilistic attribution (*Candille et al.*, 2015; *Song et al.*,  
78 2016; *Gharamti et al.*, 2017; *Mattern et al.*, 2018; *Ford*, 2019). Here, we focus our analysis on  
79 the use of ocean colour PFT observations following recent advances in data assimilation to  
80 improve marine ecosystem simulations (*Ciavatta et al.*, 2018; 2019).

81 In a probabilistic framework, the support of the joint *pdf* of observed and forecast values  
82 should be non-null in order to enable data assimilation. In other words, the prior model state  
83 and the data must be compatible with each other given their respective uncertainties in order  
84 for assimilation to be meaningful. Therefore, one important question we address in this study  
85 is the following: is the distribution of the forecast errors estimated from the ensemble (the prior  
86 distribution) compatible with the distribution of the data to be assimilated? Another important  
87 question is about how ensembles (e.g. the state of a stochastic model) can be validated – in  
88 effect, ensembles are not used solely for assimilation – other uses include: array design,  
89 probabilistic forecasting, a learning base for artificial intelligence applications, etc.

90 As a first step, we assess the empirical consistency of ensembles by means of rank histograms  
91 (also referred to as Talagrand diagrams; *Candille and Talagrand*, 2005) projecting the data

---

<sup>1</sup> The support of a probability density function (*pdf*) is the smallest closed set outside of which the *pdf* vanishes. For a *pdf* defined in  $\mathbb{R}$ , the *pdf* envelope (i.e. here the ensemble envelope) is the range between the minimum and maximum values of the support.

92 onto the model ensemble classes. Rank histograms are useful for diagnosing reliability and  
93 inferring systematic biases in an ensemble prediction system (Hamill, 2001). As a second step,  
94 we apply a *consistency diagnostic on innovations* inspired by the similarly-named diagnostic  
95 developed in Andersson (2003) and Desroziers *et al.* (2005) in the framework of the assessment  
96 of the well-posedness of data assimilation schemes. Our implementation is however specific  
97 in two ways: (1) we use that metric here outside of a data assimilation scheme, to check how  
98 consistent our ensembles are with respect to innovations; (2) we project the innovation  
99 consistency diagnostic on *Array Modes*, with the objective of permitting pattern-dependent  
100 consistency analysis. Our definition of array modes follows previous publications (Le Hénaff  
101 *et al.*, 2009; Lamouroux *et al.*, 2016; and Charria *et al.*, 2016). This is detailed below.

102 This study is organized as follows. Section 2 describes the stochastic approach to generate  
103 model ensembles, following the companion article Part 1. The specifications of the  
104 observational networks are presented in Section 3, including the description of the ocean colour  
105 PFT. A summary of the consistency analysis framework based on rank histograms and array  
106 modes is given in Section 4, which presents also a new criterion in “array space”. The results  
107 and conclusions are discussed in Sections 5 and 6.

## 108 2 Stochastic modelling

109 We used the NEMO platform (Nucleus for European Modelling of the Ocean;  
110 <http://www.nemo-ocean.eu/>; Madec, 2012) and its biogeochemical component PISCESv2  
111 (Pelagic Interactions Scheme for Carbon and Ecosystem Studies volume 2; Aumont *et al.*,  
112 2015). The ocean model domain encompasses the Bay of Biscay and the western part of the  
113 English Channel (Quattrocchi *et al.*, 2014; Vervatis *et al.*, 2016). The physical model is  
114 coupled online (one-way with high coupling frequency for the conservation of tracers) with a  
115 biogeochemical model at  $1/36^\circ$  horizontal resolution.

116 Three seasonal-range ensembles were carried out using stochastic modelling of ocean physics  
117 and biogeochemistry as part of the project SCRUM (Stochastic Coastal/Regional Uncertainty  
118 Modelling; cf. companion article Part 1 by Vervatis *et al.* 2021). Perturbations were modelled  
119 using first-order auto-regressive processes - AR(1) - in the context of stochastic perturbed  
120 parameterized tendencies (SPPT; Buizza *et al.*, 1999) and stochastic perturbed parameters  
121 (SPP; Ollinaho *et al.*, 2017). The AR(1) processes were different for each perturbed tendency  
122 or parameter, and varied under the assumption of spatiotemporal correlated scales. A  
123 deterministic free run was performed from July 2011 to November 2011, to serve as a five-  
124 month model spin-up starting from ocean analyses for the whole state vector, and then extended  
125 from December 2011 to June 2012, to serve as reference for the main physical and  
126 biogeochemical processes in the region and deduce the statistical properties for the AR(1)  
127 stochastic parameterizations. The ensemble simulations were initialized from the five-month  
128 spin-up and carried out from December 2011 to June 2012. The 40-member ensembles were  
129 used to estimate the forecast error covariance matrix perturbing different sources of model  
130 errors.

131 In Table 1, we summarize the stochastic protocol for the three ensembles, hereafter referring  
132 to them as: EnsP - perturbing physics under the assumption of atmospheric forcing  
133 uncertainties and model uncertainties in physical parameterizations; EnsB - perturbing  
134 biogeochemistry under the assumption of model uncertainties in biogeochemical sources and  
135 sinks; and EnsPB - perturbing both physics and biogeochemistry. The biogeochemical model  
136 parameters were not perturbed. The initial conditions were also not perturbed. The EnsPB is  
137 statistically identical to EnsP for model errors in ocean physics, because there is no feedback  
138 to the physics from the biogeochemical model. By contrast, EnsPB is statistically different to  
139 EnsB and EnsP because uncertainties in the physical forcing are found to have a large impact

140 on biogeochemical model properties. Also, for the EnsB the physics and hence the circulation  
141 are the same for all members. For details of the physical-biogeochemical model and for model  
142 data ensemble-based misfits, the reader is referred to the companion article Part 1.

143 The model uncertainties in physical and biogeochemical properties in the upper ocean are  
144 described in Part 1. Here we illustrate the typical patterns that we obtained for the spread in  
145 SST, sea surface height (SSH) and chlorophyll abundance at specific dates that will also serve  
146 as a reference for the analyses in Section 5. Figure 1a shows, the Bay of Biscay SST from the  
147 deterministic simulation on May 31, 2012. Fig. 1b shows the EnsP model ensemble spread for  
148 SST on the same date. Figs. 1a-b show that the ocean model solution and its associated error  
149 structures for SST are statistically consistent. In more detail, error regimes are linked to  
150 physical processes controlling SST, for example, the extension of the Ushant thermal tidal front  
151 west of Brittany in the English Channel. At that time of the year, the stratification is strong and  
152 the freshwater front (due to the Loire and Gironde rivers runoff in the Armorican shelf) leads  
153 to a large spread locally exceeding 0.8 °C. Figs. 1c-d give an insight of the model circulation  
154 patterns on February 25, 2012, depicted by SSH and their error structures. After a two-month  
155 spin-up period, the maxima of SSH model spread are collocated with mesoscale eddies  
156 observed in the abyssal plain, and with (most likely) inertial barotropic waves observed over  
157 the shelves. The perturbation mechanism works throughout the whole period. Therefore, the  
158 SSH spread continues to grow in the abyssal plain due to mesoscale decorrelation of eddies  
159 across members (not shown).

160 Figure 2a shows the total chlorophyll abundance of the deterministic simulation at the onset of  
161 the spring bloom on March 28, 2012. In Figs. 2b-d, and for the same date, we show model  
162 uncertainties in total chlorophyll for the three ensembles EnsP, EnsB and EnsPB. We  
163 emphasize two important findings: first, the model errors in physics have a larger impact on  
164 chlorophyll model uncertainties compared with those generated by perturbing the  
165 biogeochemical model source and sink terms; second, the chlorophyll spread is increased when  
166 model physics and biogeochemistry are perturbed simultaneously. However, there are also  
167 exceptions of decreasing spread (e.g. occurring locally in the presence of coherent eddies).  
168 Interestingly, an uneven ensemble spread is also observed between nanophytoplankton and  
169 diatoms chlorophyll (Figs. 2e-f). Model errors for each phytoplankton type follow the different  
170 patterns of their chlorophyll concentrations, with large uncertainties observed in the abyssal  
171 plain for nanophytoplankton, and on the shelves and the English Channel for diatoms.

### 172 **3 Observational networks**

173 In this study, although we examine the same broad categories of observations as in the  
174 companion article Part 1, we have included more satellite products observing the upper-ocean  
175 physical and biogeochemical properties. Therefore, for the sake of clarity, we have included a  
176 full description of observations. We consider global and regional satellite products monitoring  
177 the Bay of Biscay, using high-level L3 and L4 merged data from multi-missions for both real-  
178 time and delayed modes. All datasets are provided at daily frequency. This level of data quality  
179 is often required in ocean forecasting systems for model validation and data assimilation. Table  
180 2 summarizes the datasets used for the empirical consistency analysis and their specifications  
181 are given briefly in the paragraphs below.

#### 182 **3.1 SST gridded observations**

183 We chose two high-resolution SST L4 products, namely the OSTIA SST global dataset and the  
184 regional Atlantic European north west shelf (NWS) SST. High-resolution SST data are  
185 necessary when it comes to validate eddy-resolving model ensembles. Both networks are able

186 to resolve mesoscale eddy variability and frontal activity near the coasts, ranging from a few  
187 kilometres over the shelves to scales of tens of kilometres in the open ocean.

188 The OSTIA SST is provided on a global regular grid at  $0.05^\circ$  resolution (*Donlon et al.*, 2012)  
189 and is by construction free of diurnal variability with a reference depth of 10 meters, referred  
190 also as foundation SST. The reconstructed fields consist of L4 daily gap-free maps. The product  
191 includes multi-sensor measurements from both infra-red and microwave instruments, as well  
192 as in-situ observations retrieved from ships of opportunity, drifters and buoys. The data were  
193 reprocessed using optimal interpolation and include a bias correction.

194 The NWS SST product is delivered at  $0.04^\circ$  high horizontal resolution and consists of L4 daily  
195 gap-free regional maps. Calculations are derived from infra-red measurements using only  
196 night-time observations, and thus the NWS SST data are not affected by diurnal warming.

### 197 **3.2 SLA gridded and along-track observations**

198 We selected two products based on different levels of processing. The L3 along-track sea-level  
199 product is based on several altimetric missions, specifically on Envisat and Cryosat-2 satellites.  
200 The L4 gridded product was generated by reprocessing multi-mission altimetry data over many  
201 years.

202 For the along-track data, we used the filtered sub-sampled L3 product with 14 km distance  
203 between successive points along the altimetry track. The effective resolution of the data is  
204 coarser than the aforementioned distance (*Pujol et al.*, 2016). The inverse barometer response  
205 and the tides have been removed from the data (and the model, cf. companion article Part 1).

206 The L4 product is a merge of multi-mission satellites spanning the last two decades. The result  
207 is a homogenous and consistent gridded dataset at a resolution of  $0.25^\circ$ . The scales resolved  
208 by the L4 gridded product depend on various factors, such as the quality of the L3 input data,  
209 the sampling of the altimeter constellation and the optimal interpolation, which limits the  
210 capability of the data to resolve part of the mesoscale (*Chelton et al.*, 2011).

### 211 **3.3 Phytoplankton functional types ocean colour observations**

212 We selected the regionally-tuned PFT satellite model of *Brewin et al.* (2017). The approach is  
213 a modification of the three-component model of *Brewin et al.* (2010; 2015) and produces  
214 estimates of the chlorophyll concentration of four phytoplankton groups (picophytoplankton,  
215 nanophytoplankton, diatoms and dinoflagellates), as a function of the total chlorophyll  
216 concentration, accounting for the influence of SST on model parameters. The model was tuned  
217 and validated using a large dataset collected in the North Atlantic, inclusive of our region of  
218 interest. The approach was run using ocean colour data from the ESA ocean colour climate  
219 change initiative (OC-CCI) and OSTIA SST data. Products were produced on a regular grid at  
220  $0.05^\circ$  resolution. The data has been used successfully in PFT data assimilation exercises in the  
221 North Atlantic (*Ciavatta et al.*, 2018; *Skakala et al.*, 2018). The products were developed in  
222 the framework of the Copernicus Marine Environment Monitoring Service (CMEMS)  
223 “Towards Operational Size-class Chlorophyll Assimilation (TOSCA)” Service Evolution  
224 project.

225 The mismatch between biogeochemical model variables and ocean colour products has been  
226 highlighted recently in an IOCCG Report No. 19 (*IOCCG*, 2020). The biogeochemical model  
227 PISCES distinguishes two classes of chlorophyll, namely the nanophytoplankton and diatoms.  
228 These two chlorophyll classes were given this name, as the traits of these modelled functional  
229 types are closest to these two abundant groups in the natural world. However, these two classes  
230 are designed to be representative of the entire phytoplankton community, and therefore are far  
231 broader than the names imply. Each of the two model compartments encompasses a large range

232 of cell sizes, approximately representing smaller and larger phytoplankton types in the  
233 biogeochemical model. In order to compare the observations with the model ensembles, we  
234 decided to follow a size class-based categorization (*Sieburth et al.*, 1978) and group the four  
235 ocean colour PFT into two phytoplankton types and keep the total chlorophyll concentration  
236 in the data unchanged (i.e. the sum of all four types). For this, we combined the concentrations  
237 of ocean colour pico- and nano- PFT into one small class comparable with the PISCES  
238 nanophytoplankton, and the diatoms and dinoflagellates ocean colour PFT into one large class  
239 comparable with the PISCES diatoms. This categorization took into consideration many  
240 processes and biogeochemical cycles in the model, for example, not just the silicate cycle but  
241 also the carbon and nitrogen cycles, and the size-selective feeding by microzooplankton and  
242 mesozooplankton. However, it was unavoidable to combine one PFT with characteristic  
243 silicate limitation with one without (i.e. diatoms and dinoflagellates; in supplementary material  
244 we provide also an analysis without this combination). The choice to include all four PFT was  
245 supported by the fact that their sum ensures the total biomass (chlorophyll) and therefore,  
246 model and satellite total chlorophyll data can be compared like-for-like. In addition, the total  
247 chlorophyll concentration (sum of the four PFT) in the model is of the same magnitude as the  
248 total chlorophyll of other ocean colour products used for the validation of our ensembles, as  
249 shown in the companion article Part 1 and in the supplementary material.

### 250 3.4 Observational errors

251 An estimate of the observational error is given by the measurement and representativity errors  
252 (usually unknown) (*Desroziers et al.*, 2005; *Oke and Sakov*, 2008; *Janjić et al.*, 2018). The  
253 measurement errors are usually small and refer to the instrument's sensor accuracy and data  
254 processing. The representativity errors are usually large and can account for different model  
255 and observation sampling schemes, or for the physical-biogeochemical signal that is contained  
256 in the observations but is not represented in the model. Some of the data used here were  
257 retrieved from archives and included a spatial distribution of their errors. While this  
258 information is important for a dynamically heterogeneous system such as the Bay of Biscay,  
259 we chose not to use it because it would be hard to interpret pattern-dependent consistency  
260 results. Instead, we used a representative constant error for each observational network  
261 provided by the CMEMS infrastructure and the data providers.

262 Assuming that errors are uncorrelated and that the innovation variance is therefore close to the  
263 sum of model and data error variances, we considered an error standard deviation for SST equal  
264 to  $0.5\text{ }^{\circ}\text{C}$ , for SLA equal to  $0.05\text{ m}$ , and for PFT equal to  $0.3\text{ mg/m}^3$  (Table 2). The satellite  
265 chlorophyll *a* data are provided with a scaled (%) observational error in comparison to the  
266 signal and therefore, the error has spatial distribution. We chose this constant PFT error over  
267 an error that scales with the satellite chlorophyll *a* signal for the reason explained in the  
268 paragraph above. The static PFT error is representative for the region, though moderately larger  
269 for total chlorophyll and nanophytoplankton in the open ocean, and underestimated for diatoms  
270 over the shelves, if compared with the PFT errors estimated by *Brewin et al.* (2017).

271 In the following Section 4, and for each of the consistency analysis methods, we discuss the  
272 perturbation of observations to generate data distributions and the anamorphosis functions  
273 applied to transform those distributions in the proper space (*Simon and Bertino*, 2009).

## 274 4 Consistency analysis framework

275 In this study, we implemented two empirical consistency analysis methods to evaluate model  
276 ensembles with respect to observations. The first approach was based on rank histograms and  
277 the second on array modes. The calculations were performed with a toolbox (nicknamed  
278 “scrumcat”: SCRUM consistency analysis toolbox) built upon the Sequoia Data Assimilation

279 Platform (SDAP; De Mey-Frémaux, 2020: <https://sourceforge.net/projects/sequoia-dap/>). The  
280 implementation details are given below.

#### 281 4.1 Rank histograms

282 The rank histogram (Talagrand diagram; *Candille and Talagrand, 2005*) is a verification metric  
283 testing the empirical consistency between model samples and observations. An ensemble of  $m$   
284 members defines  $m + 1$  ranks including the rank of one observation sorted within the model  
285 ensemble. The steps to calculate rank histograms are (1) project ensemble samples in data space  
286 using an observation operator, (2) sorting the ensemble samples by value, (3) ranking the  
287 observed value within the sorted ensemble at each observation location, (4) tally over many  
288 observations in space and time. We note that for the rank histograms the observations are not  
289 perturbed and the model ensemble distributions are not checked for gaussianity.

290 The rank histogram is useful for determining the reliability of ensemble forecasts (*Hamill,*  
291 *2001*) and is often used as a tool to infer systematic biases of an ensemble prediction system,  
292 as well as differing spreads between model and data (e.g. small spread translates into most data  
293 falling into the outer classes of the rank histogram). The rank histogram flatness is a measure  
294 of the reliability of the ensemble prediction system and is estimated according to the definition  
295 of statistical consistency by *Anderson (1996)* and *Candille and Talagrand (2005)*:

$$296 \Delta = \sum_{i=1}^{m+1} \left( S_i - \frac{p}{m+1} \right)^2, \Delta_o = \frac{p*m}{m+1}, \delta = \frac{\Delta}{\Delta_o} \sim 1, \quad (1)$$

297 where  $p$  is the number of observations,  $m$  the number of ensemble members,  $S_i$  the outcomes  
298 in each rank (i.e. number of observations of rank  $i$ ),  $\Delta$  the deviation of the histogram from  
299 flatness,  $\Delta_o$  the expectation for a reliable system and  $\delta$  the ratio measuring the reliability for a  
300 scalar variable of the ensemble prediction system.

#### 301 4.2 Ensemble consistency in array space

302 Rank histograms provide a way to check ensemble empirical consistency in a global,  
303 distribution-like manner. In a complementary manner, we wish to complete the global analysis  
304 with pattern-selective consistency analysis, using patterns which would be (A) hierarchized,  
305 (B) representative of error covariances (as estimated from the ensemble) and (C) observable  
306 by the data array. To that end, we chose to assess consistency in the space of *Array Modes*,  
307 hereafter nicknamed “array space”.

308 We use whenever possible the unified notations of data assimilation as in *Ide et al. (1997)*. We  
309 define *array modes*  $\boldsymbol{\mu}$  as the eigenmodes of the *scaled representer matrix*  $\boldsymbol{\chi}$  defined as:

$$310 \boldsymbol{\chi} \equiv \mathbf{R}^{-\frac{1}{2}} \mathbf{H} \mathbf{P}^f \mathbf{H}^T \mathbf{R}^{-\frac{1}{2}} = \boldsymbol{\mu} \boldsymbol{\sigma} \boldsymbol{\mu}^T \quad (2)$$

311 where  $\mathbf{P}^f$  is the *forecast (prior) error covariance matrix*, here approximated as an ensemble  
312 covariance,  $\mathbf{R}$  is the *observational error covariance matrix*, and  $\mathbf{H}$  is the “classic” linear  
313 *observation operator*. Because an observing array will not be concomitant, all operators are  
314 four-dimensional in that they span time in addition to space.

315 These array modes have been discussed in several articles in the literature. First, *Le Hénaff et*  
316 *al. (2009)* presented the theoretical background for the representer matrix spectra methodology,  
317 aiming at assessing the performance of observational networks at detecting model errors. In  
318 their study, prior model errors were generated via stochastic modelling of the wind forcing and  
319 various altimetry and in-situ array deployment strategies were tested. Two other studies  
320 followed, based on array modes, by *Lamouroux et al. (2016)* and *Charria et al. (2016)*, where  
321 authors designed optimal observation network experiments based on array modes, for the  
322 future implementation of efficient integrated ocean observing systems monitoring the coastal

323 environment. Array modes are the eigenmodes of the representer matrix as approximated as an  
 324 ensemble covariance (e.g. from Section 2 ensembles). It follows that array modes meet the  
 325 pattern definition criteria (A-C) above.

326 The matrix  $\chi$  in Eq. (2) expresses how the array “sees” the model uncertainties; the diagonal  
 327 matrix  $\sigma = \text{diag}\{\sigma_k\}$  offers the same information as  $\chi$  in array space, but in a diagonal form:  
 328 it is the *array mode spectrum*. Intuitively, the higher the eigenvalues in  $\sigma$ , the better the array  
 329 can detect (and help correct) errors of a prior estimate. *Le Hénaff et al.* (2009) show that a  
 330 useful choice of a discriminating spectral value is 1. In order to clarify this choice, let us  
 331 consider the *innovation vector*  $\mathbf{d}$  and its second-order statistics:

$$332 \quad \mathbf{d} \equiv \mathbf{y}^o - \mathbf{H}\mathbf{x}^f \quad (3)$$

$$333 \quad \langle \mathbf{d}\mathbf{d}^T \rangle = \mathbf{R} + \mathbf{H}\mathbf{P}^f\mathbf{H}^T \quad (4)$$

334 where  $\mathbf{y}^o$  is the *observation vector* and  $\mathbf{x}^f$  is the *forecast (prior) state vector*. Equation (4) is  
 335 the same as Eq. (1) in *Desroziers et al.* (2005). An intuitive criterion of array performance is  
 336 as follows:

337 1) If observational errors  $\mathbf{R}$  dominate in Eq. (4), then most of the model-data discrepancies in  
 338  $\mathbf{d}$  are attributable to observational error, and observations are not being very useful at detecting  
 339 model uncertainties.

340 2) If prior state errors  $\mathbf{H}\mathbf{P}^f\mathbf{H}^T$  (the *representer matrix*) dominate, then most of the  
 341 discrepancies are attributable to model uncertainties, and observations can be expected to be  
 342 useful at identifying and correcting them.

343 Let us examine this criterion in array space. Using the orthogonality of array modes, it is easy  
 344 to show that the innovation covariance in Eq. (4) projected into array space, becomes:

$$345 \quad \mathbf{E}^i \equiv \boldsymbol{\mu}^T \mathbf{R}^{-\frac{1}{2}} \langle \mathbf{d}\mathbf{d}^T \rangle \mathbf{R}^{-\frac{1}{2}} \boldsymbol{\mu} = \mathbf{I} + \boldsymbol{\sigma}, \quad (5)$$

346 hence the choice of 1 as a discriminating spectral value. It follows:

347 *Criterion ArMI: Array performance is measured by the number of ranks  $k$  for which:*

$$348 \quad 1 \leq \sigma_k$$

349 Note that (1) this first criterion is only based on the space-time sampling scheme and on the  
 350 ensemble covariance; it does not require the actual values of observations; (2) it is not yet  
 351 expressed as an empirical consistency criterion for our ensemble. Now, when we have the  
 352 values of observations, Eq. (4) can be used to derive an empirical ensemble consistency  
 353 criterion. Innovation Eq. (3) can be formed, and  $\langle \mathbf{d}\mathbf{d}^T \rangle$  can be calculated empirically and  
 354 compared to the right-hand side of Eq. (4). In effect, the innovation spread is the result of prior  
 355 uncertainties of both the model and observations; therefore, that spread should be statistically  
 356 consistent with the sum of prior model uncertainty estimates and observational uncertainty  
 357 estimates.

358 As above, this criterion is more advantageously examined in array space, i.e. in the form of Eq.  
 359 (5). Since  $\mathbf{I} + \boldsymbol{\sigma}$  is diagonal, we focus on the diagonal of  $\mathbf{E}^i = \text{diag}\{e_k^i\}$ , i.e. the innovation  
 360 variance in array space. Our second criterion writes:

361 *Criterion ArMCA1: Empirical consistency is measured by the number of ranks  $k$  for which:*

$$362 \quad 1 - \tau \leq \frac{1 + \sigma_k}{e_k^i} \leq 1 + \tau$$

363 *with  $\tau$  the user-selected tolerance (e.g.  $\tau = 0.1$  for 10% tolerance).*

364 In this study, we use the ArM1 and ArMCA1 criteria as implemented in the ArM tools library  
 365 (*De Mey-Frémaux, pers. comm., 2020*), distributed as open source within the SDAP  
 366 assimilation platform (*De Mey-Frémaux, 2020*: <https://sourceforge.net/projects/sequoia-dap/>).  
 367 We check consistency in array space, also qualifying consistency results with their associated  
 368 multivariate patterns in observation space. In short, criterion ArM1 is an “array performance”  
 369 criterion based on how the observational array “observes” model uncertainties; criterion  
 370 ArMCA1 tests ensemble empirical consistency in terms of variances; other criteria are present  
 371 in the ArM tools but are not used here.

372 Our approach can be seen as an extension of the *Desroziers et al. (2005)* consistency diagnostic  
 373 on innovations (their Eq. (1)), which we project on the space of array modes.

374 We remove ensemble averages prior to analysis. In contrast to rank histogram analysis,  
 375 observations are perturbed, using a Gaussian random number to generate data distributions  
 376 with proper error estimates for each network (Table 2). Although the methodology allows for  
 377 nondiagonal *observational error covariance matrix*, we consider it to be diagonal in this study,  
 378 meaning that observational errors are considered statistically independent from each other. In  
 379 case of non-Gaussian distributions, an anamorphosis transformation is applied so as to  
 380 calculate array modes in the transformed space. For instance, in this study, model and data  
 381 chlorophyll samples are log-transformed prior to computing array modes.

## 382 **5 Results**

### 383 **5.1 Rank histograms consistency analysis**

384 In Figs. 3-5, we show Hovmöller plots of rank histograms as a means to analyse space-time  
 385 ensemble consistency. Under-dispersive rank histograms are characterized by a U-shape  
 386 diagram, since observations fall mostly in the outer ranks of the ensemble. In the context of a  
 387 Hovmöller plot of rank histograms, under-dispersion is illustrated by contrasting colours  
 388 between outer and central ranks (Fig. 3). In the same line of thinking, bias is depicted by having  
 389 large values only on one side of the histogram and flatness (i.e.  $\delta \sim 1$ ) is verified for evenly  
 390 distributed *pdf* values at  $\frac{1}{41} \sim 0.024$  across all ranks. Rank histograms for SST verify that EnsP  
 391 is generally under-dispersive and biased, though results vary depending on seasons and on the  
 392 use of different networks (Fig. 3). On the other hand, there are cases where observations are  
 393 evenly distributed within the ensemble spread despite persistent under-dispersion. The rank  
 394 histograms are in close agreement with the consistency analysis based on the OSTIA SST  
 395 innovation samples, presented in the companion article Part 1.

396 In order to illustrate dependency based on geographical region, we focus on three distinct areas  
 397 in the Bay of Biscay, namely the abyssal plain, the Armorican shelf and the English Channel  
 398 (Fig. 3). There is evidence of local error regimes where EnsP is either warm or cold biased  
 399 with respect to observations. From the latter we can link SST observational biases and  
 400 underestimated model errors potential to local dynamics. For instance, rank histogram flatness  
 401 is occasionally verified in the abyssal plain, whereas in the Armorican shelf and the English  
 402 Channel the ensemble is shifting from being warm biased in winter to cold biased in spring  
 403 with respect to observations. The results are associated with the spring shoaling of the  
 404 thermocline in the open ocean, as well as to coastal processes controlled by frontal activity of

405 freshwater river discharges over the Armorican shelf, and to tidal mixing in the English  
406 Channel. There is also a marked difference between the two SST networks when used to  
407 validate consistency against EnsP. The OSTIA SST fits better within the ensemble classes  
408 during winter, whilst the consistency for the NWS SST is improved during spring. Overall, the  
409 EnsP is notably under-dispersive and on average warm biased with respect to observations,  
410 whilst cold biased in the English Channel.

411 Figure 4a shows a Hovmöller plot of rank histograms and Fig. 4b a map of ranks as a means to  
412 analyse space-time ensemble consistency for sea-level. The SLA model equivalent is  
413 calculated using a mean dynamic topography as in the companion article Part 1. The gridded  
414 SLA product appears to be under-dispersive having a bias against the ensemble shifted between  
415 seasons (Fig. 4a). This bias is most likely attributed to the seasonal steric cycle contributing to  
416 sea-level variability being weak in the observations. As a consequence, rank histogram flatness  
417 is not verified throughout the series.

418 Similar results are presented in the map of rank histograms for the along-track SLA product  
419 (Fig. 4b). As in the companion article Part 1, we focus our analysis in late-winter for a few  
420 consecutive days after having a model ensemble statistical spin-up for EnsP (Fig. 1d) and also  
421 observing noticeable sea-level variability in the satellite observations. Strong biases are  
422 apparent, both in the deeper areas of the domain and in the shelves. These biases also coincide  
423 with the location of inconsistent tidal residual signals (both in the model and data) in the  
424 English Channel. Contrasting colours in the Celtic Sea and near the shelf break hint at the  
425 presence of high-frequency error processes currently unaccounted for in the model ensemble.  
426 Under-dispersion and missing errors in the high-frequency band (e.g. open boundary  
427 conditions in the shelves, which do not deal properly with high-frequency processes), is also  
428 confirmed by the fact that there are only a few central ranks depicted in the map.

429 In Fig. 5, we present Hovmöller plots of rank histograms as a means to analyse ensemble  
430 consistency with respect to different classes of chlorophyll. For this, we use the PFT dataset  
431 against the model ensemble EnsPB which has the largest spread among the three ensembles  
432 (Figs. 2b-d). Rank histograms for EnsPB verify that all ensembles are under-dispersive, with  
433 EnsB being the least dispersive and EnsPB the most dispersive. The chlorophyll rank  
434 histograms for EnsP (not shown) are comparable to those of EnsPB (Fig. 5), but with lower  
435 levels of consistency. Rank histograms indicate a rather long statistical biogeochemical spin-  
436 up period, on the order of 3 months (verified by rank histogram flatness for evenly distributed  
437 *pdf* values), and persistent model underestimation of chlorophyll abundance. The spin-up  
438 period corresponds only to the winter season characterised by low primary production. After  
439 the spin-up period, the consistency is in general improved for total chlorophyll, as a result of  
440 the onset of the nanophytoplankton spring bloom. On the other hand, diatoms appear to  
441 degenerate the consistency of total chlorophyll. During successive spring blooms,  
442 inconsistencies are seen in the peak, but better agreement in the onset and relaxation of the  
443 blooms. This is explained by the fact that ocean colour peak values are of an order of magnitude  
444 larger compared with model outputs, whilst PFT data fall within the ensemble spread during  
445 the onset and relaxation phases of a bloom event. Overall, the most consistent configurations  
446 are the nanophytoplankton in early-spring (i.e. primary bloom mainly in the abyssal plain) and  
447 the diatoms in late-spring (i.e. secondary bloom mainly coastal and over the shelves).

## 448 5.2 Array-space empirical consistency analysis

449 In Fig. 6a, we show in a Hovmöller plot the variations in time of OSTIA SST array mode  
450 spectra vs. modal rank (i.e. the rank of the array mode), including an EnsP consistency check  
451 against OSTIA SST data. We used 39 array modes, corresponding to the 39 “degrees of  
452 freedom” characterizing a 40-member de-biased ensemble. The ensemble spin-up period

453 appears as a period of low eigenvalues near the beginning of the time series, lasting from a few  
454 days (in the head of the spectra) to a few weeks (in the tails). The array performance at  
455 “detecting” model uncertainties, as per criterion ArM1, appears satisfactory from spin-up time  
456 to the end of the series, with values above 1 across the spectra. Eigenvalues get larger in the  
457 last two months, likely reflecting both the onset of a seasonal thermocline and upper-ocean  
458 processes within the error subspace generated by our stochastic protocol (e.g. in response to  
459 wind uncertainties). The loss of empirical consistency of EnsP with respect to OSTIA SST as  
460 per criterion ArMCA1 appears as white “pixels”. Consistency appears to be almost always  
461 verified along the dominant array modes, mostly associated with large-scale patterns (as  
462 illustrated in Fig.7). The number of ensemble inconsistency cases increases as one moves to  
463 the tails of the spectra, mostly associated with smaller-scale and coastal patterns. Overall,  
464 pattern consistency is fairly good with OSTIA SST.

465 Figure 6b shows the impact of subsampling the OSTIA data, mimicking the impact of using a  
466 lower-resolution SST product for ensemble validation. However, although the general pattern  
467 is similar, the statistical significance of the consistency analysis degrades significantly; only  
468 ~50% of the array modes pass the ArM1 criterion (many more than with the product without  
469 subsampling) and appear inconsistent in the ArMCA1 criterion. It is important to note that this  
470 degradation is also felt for the dominant array modes, mostly associated with large-scale  
471 patterns, except possibly near the end of the time series. Such a (random and systematic)  
472 subsampling leads to a significant loss of information that could be used to constrain the high-  
473 resolution model ensemble, if EnsP was to be used as an estimate of the model errors. Also,  
474 high-resolution SST products appear necessary when it comes to validating eddy-resolving  
475 ensembles, such as  $1/36^\circ$  here.

476 Is this degradation dependent upon the data product? The NWS SST observations are of higher  
477 resolution compared with the OSTIA SST and therefore are able to “detect” a broader range of  
478 model errors, including in the tail modes (not shown). A subsampling rate of one-sixth brings  
479 up a resolution similar to the one-fifth subsampling of OSTIA dataset (~25 km). The  
480 consistency analysis results were found to be comparable for both networks (Figs. 6b, c).

481 Let us now turn to data-space patterns of array modes, i.e. the eigenvectors over which our  
482 empirical ensemble consistency analysis of EnsP is carried out. Figure 7 shows examples of  
483 such patterns for both one-fifth subsampled OSTIA (7a-c) and one-sixth NWS (7d), along with  
484 the result of the ArMCA1 criterion for each mode. We focus on May 31, 2012 which is  
485 characterized by a slightly higher energy in the spectra compared with neighbouring periods  
486 (Figs. 6b-c). Unlike rank histograms depicted in Fig. 3, the OSTIA and NWS array modes are  
487 very similar (as illustrated in Figs.7a, d for mode 1), because the subsampled observational  
488 schemes are similar, but consistency analysis results do not have to be. Not surprisingly for a  
489 regularly-spaced dataset, the first array mode patterns (Fig.7a, d) resembles the structure of the  
490 SST spread (Fig. 1b).

491 For the datasets at hand, the dominant patterns appear to be mostly typical of large-scale and  
492 mesoscale processes (Figs.7a, b); seeing zero amplitudes in low-order modes for the English  
493 Channel indicates that the array will not be able to correct the SST statistically dominant model  
494 errors. The spatial scales appear to decrease as one moves to the tail of the spectrum, and the  
495 coastal features become more pronounced (as illustrated in Fig.7c); simultaneously, the  
496 ensemble consistency is verified less frequently.

497 When an ensemble is found to be inconsistent with observations along a particular direction  
498 (here, an array mode), a plausible reason is that other error processes must be active in the  
499 model in addition to the ones which are at work across the ensemble. This can be verified by  
500 checking for model underdispersion. However, as is the case in the examples above,

501 insufficient resolution of the observations can also lead to statistical inconsistency, impacting  
502 mostly small scales and coastal scales.

503 In Fig. 8a, we show a Hovmöller plot of variations in time of the gridded SLA array mode  
504 spectra, including an EnsP consistency check against gridded SLA data. Again, we used 39  
505 array modes. The dataset is relatively coarse at  $\sim 25$  km resolution and therefore is not  
506 subsampled. As before, spectral values smaller than one (light grey in Fig.8a) are associated  
507 with poor array performance at “detecting” model uncertainties as per the ArM1 criterion, and  
508 white areas depict inconsistent array modes as per ArMCA1 criterion. On average, on the order  
509 of 20 array modes sit above the observational noise floor. However, only a fraction of those  
510 modes, on the order of 10 on average, show pattern consistency as per ArMCA1 criterion.  
511 Overall, consistency of EnsP with respect to the gridded SLA product is not as good as it was  
512 with respect to the SST product.

513 Let us turn to the patterns of eigenvectors in an effort to explain the weak consistency. In Figs.  
514 8b-e, we show examples of consistent and inconsistent patterns of gridded SLA array modes  
515 on April 30, 2012. Consistency is verified for the first two array modes shown (ranks 1 and 5),  
516 featuring large-scale sea-level gradients over the whole domain including the shallow Celtic  
517 Sea and the Eastern Biscay shelf (Fig. 8b-c). Both higher array modes (ranks 10 and 20) shown  
518 in Figs. 8d-e are inconsistent – they seem to be dominated by mesoscale and submesoscale sea-  
519 level signals (or at least, their  $0.25^\circ$  representation) as well as regional-scale sea-level signals  
520 over the shelves and in the English Channel. Small-scale inconsistent higher array modes in  
521 the abyssal plain can be explained by the mesoscale low-frequency decorrelation between  
522 members, generating phase differences between mesoscale features. From these results we can  
523 infer that assimilating the  $0.25^\circ$  gridded product in EnsP would probably lead to unsatisfactory  
524 results at the mesoscale and in coastal areas, since error vicinities in the model ensemble and  
525 in observations seem at least partly inconsistent with each other in that scale range. Of course,  
526 one can introduce representativity errors in the observational error covariance matrix to make  
527 things consistent again, but small scales are expected to be inadequately controlled by  
528 assimilation.

529 We now check the consistency of EnsP against the along-track sea-level data. We focus on the  
530 same period and on the same tracks as in the rank histograms analysis with the same dataset,  
531 starting on February 25, 2012, for 10 consecutive days. The analysis is performed considering  
532 together all tracks that would be assimilated in a 10-day assimilation cycle. Therefore, each  
533 array mode, characteristic of the 10-day multi-track network, spans both space and time. Figure  
534 9a shows the array mode spectrum over the 10-day period. As per criterion ArM1, almost all  
535 eigenvalues are larger than the observational noise threshold (except a few tail modes) and as  
536 per criterion ArMCA1 most of them are consistent. This is significantly improved compared  
537 to the consistency of EnsP with respect to the gridded sea-level product.

538 In order to get a bit further into which patterns and underlying processes appear to be  
539 (in)consistent, we examine the corresponding eigenvectors. Figures 9b-c show the first two  
540 ranks of array modes for the multi-track network; Figs. 9d-e are the same for the higher array  
541 modes with ranks 19 and 20 respectively. Both Figs. 9b and 9c show large-scale sea-level  
542 gradients, encompassing open-ocean and shelf regions. The patterns in Fig. 9b are generally of  
543 a bipolar nature with near zero amplitudes along the continental shelf break for most of the  
544 tracks in the region, hinting at shelf/open-ocean exchange processes at work in the model’s  
545 error subspace. In Figs. 9d-e, we show two examples of higher array modes, with ranks 19 and  
546 20 being consistent and inconsistent, respectively. Higher array modes show a wide range of  
547 spatial scales, including large-scale over the shelves and tracks characteristic of mesoscale  
548 signals in the open-ocean, with gradients more pronounced than in the gridded product. In most

549 cases, the higher array modes are consistent and because it is above observational noise (as  
550 shown in Fig. 9a) their consistency is meaningful. We therefore have evidence that along-track  
551 SLA data could, if assimilated using EnsP for covariances, bring useful consistent information  
552 at large scales and at the mesoscale.

553 Central to this work is the consistency of biogeochemical model ensembles against our PFT  
554 product as derived from satellite ocean colour measurements. In Fig. 10, we show Hovmöller  
555 plots of variations in time of univariate PFT array mode spectra, including EnsP, EnsB and  
556 EnsPB consistency analysis against PFT data, using again 39 array modes. Calculations were  
557 performed in log array space (i.e. involving log transformation from data space prior to array  
558 mode calculations), without subsampling data. The array mode spectra in Fig. 10 exhibit strong  
559 variations in time, with several peaks corresponding to differential blooms across ensembles  
560 for the various chlorophyll classes during the onset and relaxation of those events. We note  
561 however that this is a gridded product with gaps, so part of those variations are also explained  
562 by the spatial data coverage variations with time (Fig. 10a; superimposed black line). This does  
563 not question our statistical approach, since poor data coverage will logically lead to poorer  
564 array performance at detecting ensemble variance; but that dependency must be kept in mind  
565 when interpreting the results in terms of processes.

566 With the ensemble strategy we adopted, one can examine the array performance and the  
567 consistency individually for physical (EnsP) and biogeochemical (EnsB) uncertainties, and for  
568 both (EnsPB) together, all in the same PFT data space; in that manner, we can see the individual  
569 and combined effects of both components (Fig. 10). As already noted in Part 1 article, it appears  
570 that physical perturbations (EnsP) have greater impact on biogeochemical model errors than  
571 biogeochemical sources and sinks perturbations. The higher EnsP ensemble variance for PFT  
572 variables explains also the shorter statistical spin-up period on the order of 1 month, whereas  
573 EnsB shows a spin-up period on the order of 3 months, shown as overly inconsistent (i.e. the  
574 first three months in Figs. 10a-c vs. 10d-f). It also appears that visually “adding” the  
575 nanophytoplankton and diatoms spectra does not come close to describing the total chlorophyll  
576 spectra, in particular for EnsP and EnsPB (Figs. 10a-c and 10g-i), whilst being partially verified  
577 for EnsB (Figs. 10d-f). The spectra derived from diatoms in EnsP and EnsPB appears not to  
578 contribute to the total chlorophyll spectra, the latter having eigenvalues similar to those of  
579 nanophytoplankton.

580 We see three possible classes of explanations to this result: (1) it is likely that the constant  
581 observational error of  $0.3 \text{ mg/m}^3$  over the whole domain and for all classes, is unrealistic,  
582 particularly for diatoms exhibiting higher local uncertainties in the shelves and coastal regions;  
583 (2) it is possible that our decomposition of total chlorophyll into two size classes (or four binned  
584 into two as discussed in Section 3.3) is not entirely relevant process-wise (i.e. combining one  
585 PFT with characteristic silicate limitation with one without), effectively leading to a statistical  
586 ill-posed issue; (3) total chlorophyll as resulting from a variety of processes (especially when  
587 physical perturbations are applied) will have higher statistical complexity (higher number of  
588 “degrees of freedom”) than either nanophytoplankton or diatoms taken independently, leading  
589 to “purple” spectra with more tail modes not represented on Figs. 10a and 10g, while Figs. 10c  
590 and 10i show “redder” diatoms spectra.

591 Figure 11 shows examples of low and high gridded PFT array modes and consistency status of  
592 EnsP and EnsB vs. the PFT dataset (EnsPB eigenvectors are similar to EnsP; not shown). Let  
593 us focus on March 28, 2012, where for this specific date the PFT are gap-free over the Bay of  
594 Biscay facilitating our analysis at all scales. As for rank histogram analyses, we choose to  
595 investigate the period during early spring, because consistency can be easily explained by the  
596 onset of the chlorophyll abundance primary bloom (Fig. 2). In effect, for a positive oriented

597 variable such as chlorophyll, bloom-related model spread increases when abundance increases  
598 during a bloom event, therefore chances for pattern consistency are higher. We note that this is  
599 a very different result compared to rank histograms which show degraded consistency during  
600 blooms. This is because rank histograms are sensitive to model-data biases amplified during  
601 blooms (in contrast to array modes free from biases). The latter is also supported by the fact  
602 that the biogeochemical model ensembles are under-dispersive, due to model underestimation  
603 of chlorophyll abundance.

604 The PFT first array modes using EnsP and EnsB show consistent patterns and are compatible  
605 in structure with the model's second-order statistical moments respectively (Figs. 11a, f vs.  
606 Figs. 2b, c). Another remark is that the first array modes using EnsP and EnsB are different  
607 (Figs. 11a-c vs. Figs. 11d-f). When physics is perturbed (i.e. EnsP) the total chlorophyll  
608 eigenvectors appear as large-scale open ocean patterns controlled by nanophytoplankton  
609 primary production (Figs. 11a-b). When biogeochemical source and sink terms are perturbed  
610 (i.e. EnsB) we observe mesoscale patterns for all classes, with nanophytoplankton again  
611 defining the eigenvectors of total chlorophyll (Figs. 11d-e). For both ensembles, diatom  
612 measurements appear to detect model errors mostly on the shelves, but their contribution for  
613 the first array mode for total chlorophyll is limited (Figs. 11c, f).

614 For the higher modes we observe several combinations of consistent and inconsistent modes  
615 across phytoplankton types and total chlorophyll. In general, when both phytoplankton types  
616 are found to be consistent (inconsistent), then total chlorophyll is also likely to be found  
617 consistent (inconsistent). Mixed consistency results between chlorophyll classes are also  
618 observed (less frequently), but the non-mixed examples are the most informative to investigate  
619 (Figs. 11g-l). When physics is perturbed the PFT array mode 25 is able to detect model errors  
620 for both chlorophyll classes at small-scale (Figs. 11g-i). Surprisingly, for the higher array  
621 modes using EnsP the diatoms appear to contribute to the spatial variability of total chlorophyll  
622 at small-scale, especially in the shelves (Figs. 11g-i). By contrast, when biogeochemical source  
623 and sink terms are perturbed the model-based PFT array mode 25 appears not to be consistent  
624 with PFT data at the mesoscale (Figs. 11j-l). The fact that EnsB generates chiefly eigenvectors  
625 with mesoscale patterns, implies that having identical ocean physics across all members limits  
626 the potential of the PFT measurements to detect biogeochemical model errors at small-scale  
627 (as shown for example for EnsP in Figs. 11g-i).

## 628 **6 Discussion and conclusions**

629 This study is Part 2 of a two-part series following a companion article Part 1 aimed at  
630 generating ocean model ensembles. Part 2 article focuses on: (1) sophisticated ensemble  
631 model-data comparison methods, one of them published here for the first time based on a new  
632 criterion in the space of array modes, and (2) applying those methods for the first time with  
633 phytoplankton functional type data derived from ocean colour.

634 The empirical consistency analysis focused on satellite observations (SST, SLA and ocean  
635 colour), in concert with model ensembles of ocean physics and biogeochemistry. We used a  
636 high-resolution configuration for the Bay of Biscay, as a means to investigate the model error  
637 subspace generated by our stochastic protocol in a companion article Part 1, both in the open  
638 and coastal ocean. In order to better understand the couplings between physics and  
639 biogeochemistry we examined three model ensemble experiments: perturbing only physics,  
640 perturbing only biogeochemical source and sink terms, and perturbing both simultaneously.  
641 Below, we synthesize results from the two consistency methods in an attempt to assess the  
642 reliability of model ensembles with respect to observations.

643 Rank histograms for physical properties showed a large dependency based on geographical  
644 region and on season for both SST and SLA networks. In most cases, rank histograms showed  
645 large biases between model and data, and were limited because of model ensemble underspread  
646 and because of weak variability in the observations. In addition, there were high-frequency  
647 errors in the observations that were not present in the model ensemble. Rank histograms for  
648 the ocean colour PFT showed persistent model underestimation and underspread for  
649 chlorophyll abundance, with some improvement for nanophytoplankton after a spin-up period  
650 on the order of 3 months during the winter low primary productivity. Rank histograms have  
651 been used successfully as a reliability tool identifying on several occasions consistent model-  
652 data configurations and attributing this result to physical and biogeochemical processes, such  
653 as the spring shoaling of the thermocline, the frontal activity in the shelves explained by river  
654 plume migration, the tidal mixing in the English Channel (*Karagiorgos et al., 2020*), and to a  
655 lesser extent the chlorophyll abundance during spring, mainly on the onset and relaxation of  
656 bloom events.

657 Array modes were performed using innovation samples in array space and therefore, model-  
658 data distributions were free from biases hindering empirical consistency of error patterns (as  
659 in some cases for the rank histograms). A typical result using array modes in the context of  
660 stochastic modelling is that insufficient resolution of the verifying observations can lead to  
661 statistical inconsistency, impacting mostly small and coastal scales. High-resolution datasets  
662 (and model ensembles) can improve array modes consistency at small-scale, as long as the  
663 following assumptions are valid: (a) data errors are uncorrelated, (b) data errors are small in  
664 comparison to model-data misfit and (c) the model ensembles are not under-dispersive. When  
665 ensembles are found to be inconsistent with observations along an array mode, a plausible  
666 reason is that other error processes must be active in the model, in addition to the ones at work  
667 across the ensemble. For instance, model errors in physical processes currently unaccounted  
668 for are attributable to residual tidal errors (due to local tidal fronts and occasional Kelvin waves  
669 propagating along the coasts), to the non-isostatic response of atmospheric pressure, and to  
670 high-frequency errors in open boundaries over the shelves (*Vervatis et al., 2021*). In addition,  
671 complex processes in the biogeochemical model are simplified by making use of only a few  
672 model parameters. Biogeochemical parameterizations controlling the growth rate and grazing  
673 of phytoplankton classes in the model are based on approximations and therefore, stochastic  
674 modelling of these parameters may improve model performance (*Garnier et al., 2015*).

675 Array modes for both SST networks (i.e. OSTIA and NWS) showed that despite their  
676 differences in production and resolution, if subsampled, can provide comparable information  
677 to detect model errors, from shortly after the ensemble spin-up to the end of the series (in  
678 contrast to rank histograms). Their consistency was verified for the large-scale and open ocean,  
679 but also for the small-scale and coastal SST patterns. Array modes pattern consistency between  
680 the model ensemble perturbing physics and the sea-level datasets was in general verified at  
681 large-scale. The gridded SLA product showed unsatisfactory results at the mesoscale and in  
682 coastal areas, whereas the along-track sea-level showed that, if assimilated, can bring  
683 potentially useful information at the mesoscale.

684 The most important findings for the array modes using PFT against ocean biogeochemical  
685 ensembles can be summarised as follows. A large ensemble variance can lead to an improved  
686 PFT array performance, considering biogeochemical model uncertainties stemming mainly  
687 from stochastic physics. Consistency results can be further enhanced if we consider a two-way  
688 coupling and a feedback from the bio-optical model (due to changes in chlorophyll, leading to  
689 increased model spread) to the physical model for the solar radiation penetration in the water  
690 column. We performed a few simulations (not shown) but we decided not to activate the two-  
691 way coupling, because the onset and relaxation of the spring bloom may occur earlier and more

692 rapidly in the model compared with observations in the region (*Gutknecht et al.*, 2019).  
693 Additive spectra of nanophytoplankton and diatoms cannot describe total chlorophyll spectra.  
694 One plausible explanation is that EnsP and EnsB state vectors (i.e. their anomalies from the  
695 ensemble mean) are not statistically independent (perhaps as an artefact of the limited ensemble  
696 size). This is proven also by the fact that ensemble variance for EnsPB sometimes can be  
697 unexpectedly small, as if physics and biogeochemistry processes compensated, cf. companion  
698 article Part 1. Consequently, there is no reason for array mode eigenspectra to be additive on a  
699 rank-by-rank basis nor hierarchised in the same way (in contrast to other spectra, e.g. Fourier  
700 and Wavelet Transform). Another point is that there is no reason our ensembles should be  
701 orthogonal, because they are not built this way. The array mode spectra give an idea of the  
702 number of “degrees of freedom” of the ensembles in data space, and of the very high temporal  
703 variability of that number. For instance, the bloom periods are characterized by larger numbers  
704 of “degrees of freedom”, apparently peaking at 39 (40 minus mean). Other plausible  
705 explanations may be the unrealistic observational error set globally at  $0.3 \text{ mg/m}^3$  for all PFT  
706 and also, the total chlorophyll decomposition of four PFT binned into two not being entirely  
707 relevant process-wise.

708 In general, the data PFT errors are not negligible and are related to the size-class (*Brewin et*  
709 *al.*, 2017; *Laiolo et al.*, 2021). Errors are lower for the small size classes and higher for the  
710 large size classes (with dinoflagellates having the higher errors of all PFT), and increased with  
711 increasing chlorophyll. This is owing to the nature of how the satellite PFT error is computed.  
712 Chlorophyll *a* uncertainties (e.g. root mean square error and bias) for each PFT are assigned  
713 based on optical water type membership at a given satellite pixel (*Brewin et al.*, 2017). It is  
714 expected that observational errors are cross-correlated, both across functional types and  
715 spatially for each PFT. In fact, each PFT is computed as a function of the same total chlorophyll  
716 product, and the same procedure is adopted to compute the errors of the different PFT (*Brewin*  
717 *et al.*, 2017). The total chlorophyll product itself is expected to have spatially correlated errors.  
718 This is due to both atmospheric effects (e.g. cloud cover) and in-water sources of errors, in  
719 particular in the coastal zone, due to the coloured dissolved organic matter (CDOM) inputs that  
720 degrade ocean colour near the coast (e.g. terrestrial CDOM in river plumes). This is also  
721 confirmed by the PFT spatial error gradients and the gradients of the concentrations  
722 themselves, being stronger for diatoms than nanophytoplankton in the coastal zone (*Brewin et*  
723 *al.*, 2017; *Ciavatta et al.*, 2018). Cross-correlations are expected to be lower in the open ocean  
724 and deeper areas. The representation of such uncertainty and correlated errors has been so far  
725 neglected in ocean colour data assimilation (*IOCCG*, 2020). In our application, if we used such  
726 PFT uncertainties and data correlations, the array modes consistency would have been  
727 degraded, pertaining to the fact that ocean colour datasets are provided in high-resolution (and  
728 therefore, spatial data correlations may not be negligible).

729 Pattern-selective consistency analysis showed that low-rank eigenvectors appear as large-scale  
730 and mesoscale, mainly controlled by nanophytoplankton in the open ocean; diatom  
731 measurements appear to detect model errors mostly on the shelves. For the higher modes,  
732 biogeochemical model errors appear to be detected at small-scale only when physics is  
733 perturbed. By contrast, when only biogeochemical source and sink terms are perturbed, the  
734 model-based high-rank modes appear not to be consistent with PFT data at the mesoscale,  
735 limiting also the potential to detect model errors at small-scale due to identical physical  
736 processes across all members.

737 We recommend using methods adapting and estimating the  $R$  and  $Q$  observational and model  
738 error covariance matrices, respectively. The ArM methodology can provide pattern- and scale-  
739 dependent ensemble consistency checks, facilitating the qualification of ensembles to provide  
740 useful error covariance estimates in regional systems for ulterior coastal downscaling.

741 The next step of this study will be dedicated to the investigation of ensemble consistency based  
742 on array modes: (1) performing multivariate analysis between physics and biogeochemistry  
743 (e.g. SST vs. total chlorophyll), and between phytoplankton size classes, using information for  
744 the spatial distribution of observational errors; and (2) taking into account correlated  
745 observational errors. To that end, an additional criterion will be introduced to check the  
746 diagonality of the innovation covariance matrix in the space of array modes.

747 **Acknowledgments.** This work was carried out as part of the Copernicus Marine Environment  
748 Monitoring Service (CMEMS) “Stochastic Coastal/Regional Uncertainty Modelling  
749 (SCRUM)” Service Evolution project. CMEMS is implemented by Mercator Ocean  
750 International in the framework of a delegation agreement with the European Union. Part of this  
751 research was also made possible through the IKY Scholarships Programme and co-financed by  
752 the European Union (European Social Fund-ESF) and Greek national funds through the action  
753 entitled “Reinforcement of Postdoctoral Researchers”, in the framework of the Operational  
754 Programme “Human Resources Development Programme, Education and Lifelong Learning”  
755 of the National Strategic Reference Framework (NSRF) 2014-2020. The contribution of P. De  
756 Mey-Frémaux and N. Ayoub is supported by Centre National de la Recherche Scientifique  
757 (CNRS). We acknowledge the use of the ECMWF's computing and archive facilities in this  
758 research. This work was also supported by computational time granted from the Greek  
759 Research & Technology Network (GRNET) in the National HPC facility – ARIS – under  
760 project ID PA002007. We thank five anonymous reviewers for their constructive comments.

## 761 **References**

- 762 Anderson, J., 1996. A method for producing and evaluating probabilistic forecasts from  
763 Ensemble model integrations, *J. Climate*, 9, 1518–1530.
- 764 Andersson, E., 2003. Modelling of innovation statistics. Proceedings of the Workshop on  
765 recent developments in data assimilation for atmosphere and oceans. ECMWF, Reading,  
766 UK, 153-164.
- 767 Anderson, J.L., 2009. Spatially and temporally varying adaptive covariance inflation for  
768 ensemble filters. *Tellus A*, 61: 72-83. doi:10.1111/j.1600-0870.2008.00361.x.
- 769 Auclair, F., P. Marsaleix, and P. de Mey-Frémaux, 2003. Space-time structure and dynamics  
770 of the forecast error in a coastal circulation model of the Gulf of Lions. *Dynamics of*  
771 *Atmospheres and Oceans*, Elsevier, 36 (4), pp.309-346. 10.1016/S0377-0265(02)00068-4.
- 772 Aumont, O., Ethé, C., Tagliabue, A., Bopp, L., and Gehlen, M., 2015. PISCES-v2: an ocean  
773 biogeochemical model for carbon and ecosystem studies. *Geosci. Model Dev.* 8, 2465-2513.
- 774 Brewin, R. J. W., Ciavatta S., Sathyendranath S., Jackson T., Tilstone G., Curran K., Airs R.  
775 L., Cummings D., Brotas V., Organelli E., Dall’ Olmo G., and Raitos D. E., 2017.  
776 Uncertainty in Ocean-Colour Estimates of Chlorophyll for Phytoplankton Groups. *Front.*  
777 *Mar. Sci.* 4:104. doi: 10.3389/fmars.2017.00104.
- 778 Brewin, R. J. W., S. Sathyendranath, T. Hirata, S. J. Lavender, R. Barciela, and N. J. Hardman-  
779 Mountford, 2010. A three-component model of phytoplankton size class for the Atlantic  
780 Ocean, *Ecol. Modell.*, 221, 1472–1483, doi: 10.1016/j.ecolmodel.2010.02.014.
- 781 Brewin, R. J. W., Sathyendranath, S., Jackson, T., Barlow, R., Brotas, V., Airs, R., et al., 2015.  
782 Influence of light in the mixed layer on the parameters of a three-component model of  
783 phytoplankton size structure. *Remote Sens. Environ.* 168, 437–450. doi:  
784 10.1016/j.rse.2015.07.004.

- 785 Candille, G., and Talagrand, O., 2005. Evaluation of probabilistic prediction systems for a  
786 scalar variable. *Q. J. R. Meteorol. Soc.*, 131, 2131-2150.
- 787 Candille, G., Brankart, J.-M., and Brasseur, P., 2015. Assessment of an ensemble system that  
788 assimilates Jason-1/Envisat altimeter data in a probabilistic model of the North Atlantic  
789 ocean circulation. *Ocean Science* 11, 425–438.
- 790 Charria, G., Lamouroux, J., and De Mey, P., 2016. Optimizing observational networks  
791 combining gliders, moored buoys and FerryBox in the Bay of Biscay and English Channel.  
792 *Journal of Marine Systems*, 162, 112-125, <http://dx.doi.org/10.1016/j.jmarsys.2016.04.003>.
- 793 Chelton, D. B., Schlax, M. G., Samelson, R. M., 2011. Global observations of nonlinear  
794 mesoscale eddies, *Prog. Oceanogr.*, 91, 167–216, doi:10.1016/j.pocean.2011.01.002.
- 795 Ciavatta, S., Brewin, R. J. W., Skákala, J., Polimene, L., de Mora, L., Artioli, Y., and Allen, J.  
796 I., 2018. Assimilation of Ocean-Colour Plankton Functional Types to Improve Marine  
797 Ecosystem Simulations. *Journal of Geophysical Research: Oceans*, 123(2), 834–854,  
798 <https://doi.org/10.1002/2017JC013490>.
- 799 Ciavatta, S., Kay, S., Brewin, R. J. W., Cox, R., Di Cicco, A., Nencioli, F., and Tsapakis, M.,  
800 2019. Ecoregions in the Mediterranean Sea Through the Reanalysis of Phytoplankton  
801 Functional Types and Carbon Fluxes. *Journal of Geophysical Research: Oceans*, 124(10),  
802 6737–6759, <https://doi.org/10.1029/2019JC015128>.
- 803 Desroziers, G., Berre, L., Chapnik, B., Poli, P., 2005. Diagnosis of observation, background  
804 and analysis-error statistics in observation space. *Q. J. R. Meteorol. Soc.* 131, 3385–3396.  
805 <https://doi.org/10.1256/qj.05.108>.
- 806 Donlon, C.J., Martin, M., Stark, J., Roberts-Jones, J., Fiedler, E., and Wimmer, W., 2012. The  
807 Operational Sea Surface Temperature and Sea Ice Analysis (OSTIA) system. *Remote  
808 Sensing of the Environment*, doi: 10.1016/j.rse.2010.10.017 2011.
- 809 Evensen, G., 2003. The Ensemble Kalman Filter: theoretical formulation and practical  
810 implementation. *Ocean Dyn.* 53, 343-367.
- 811 Ford, D., 2019. Assessing the role and consistency of satellite observation products in global  
812 physical-biogeochemical ocean reanalysis. *Ocean Sci. Discuss.* 1–25.  
813 <https://doi.org/10.5194/os-2019-118>.
- 814 Garnier, F., Brankart, J.M., Brasseur, P., and Cosme, E., 2016. Stochastic parameterizations of  
815 biogeochemical uncertainties in a 1/4° NEMO/PISCES model for probabilistic comparisons  
816 with ocean color data. *Journal of Marine Systems*, 155, 59–72.
- 817 Ghantous, M., Ayoub, N., De Mey-Frémaux, P., Vervatis, V., and Marsaleix, P., 2020.  
818 Ensemble downscaling of a regional ocean model. *Ocean Modell.*, 145,  
819 <https://doi.org/10.1016/j.ocemod.2019.101511>.
- 820 Gharamti, M. E., J. Tjiputra, I. Bethke, A. Samuelsen, I. Skjelvan, M. Bentsen, and L. Bertino,  
821 2017. Ensemble data assimilation for ocean biogeochemical state and parameter estimation  
822 at different sites, *Ocean Modelling*, 112, 65-89,  
823 <https://doi.org/10.1016/j.ocemod.2017.02.006>.
- 824 Gutknecht, E., Reffray, G., Mignot, A., Dabrowski, T., and Sotillo, M. G., 2019. Modelling the  
825 marine ecosystem of Iberia–Biscay–Ireland (IBI) European waters for CMEMS operational  
826 applications, *Ocean Sci.*, 15, 1489–1516, <https://doi.org/10.5194/os-15-1489-2019>.

- 827 Hamill, T.M., 2001. Interpretation of Rank Histograms for Verifying Ensemble Forecasts.  
828 Mon. Wea. Rev., 129, 550-560, [https://doi.org/10.1175/1520-0493\(2001\)129<0550:](https://doi.org/10.1175/1520-0493(2001)129<0550:)  
829 IORHFV>2.0.CO;2.
- 830 Ide, K., Courtier P., Ghil M., and Lorenc A., 1997. Unified notations for data assimilation:  
831 operational, sequential and variational. *J. Met. Soc. Japan*, 75(1B):181-189.
- 832 IOCCG, 2020. Synergy between Ocean Colour and Biogeochemical/Ecosystem Models.  
833 Dutkiewicz, S. (ed.), IOCCG Report Series, No. 19, International Ocean Colour  
834 Coordinating Group, Dartmouth, Canada.
- 835 Janjić, T., Bormann N., Bocquet M., Carton J.A., Cohn S.E., Dance S.L., Losa S.N., Nichols  
836 N.K., Potthast R., Waller J.A., Weston P., 2018. On the representation error in data  
837 assimilation, *Q J R Meteorol Soc.*,144,1257–1278.
- 838 Karagiorgos, J., V. Vervatis, and S. Sofianos, 2020. The Impact of Tides on the Bay of Biscay  
839 Dynamics. *J. Mar. Sci. Eng.*, 8, 617.
- 840 Kourafalou V.H., P. De Mey, M. Le Hénaff, G. Charria, C.A. Edwards, R. He, M. Herzfeld,  
841 A. Pasqual, E.V. Stanev, J. Tintoré, N. Usui, A.J. Van Der Westhuysen, J. Wilkin and Zhu,  
842 X., 2015a. Coastal Ocean Forecasting: system integration and validation. *Journal of*  
843 *Operational Oceanography*, <http://dx.doi.org/10.1080/1755876X.2015.1022336>.
- 844 Kourafalou, V.H., P. De Mey, J. Staneva, N. Ayoub, A. Barth, Y. Chao, M. Cirano, J. Fiechter,  
845 M. Herzfeld, A. Kurapov, A.M. Moore, P. Oddo, J. Pullen, A.J. van der Westhuysen and  
846 Weisberg, R.H., 2015b. Coastal Ocean Forecasting: science foundation and user benefits.  
847 *Journal of Operational Oceanography*, <http://dx.doi.org/10.1080/1755876X.2015.1022348>.
- 848 Lamouroux, J., G. Charria, P. De Mey, S. Raynaud, C. Heyraud, P. Craneguy, F. Dumas and  
849 Le Hénaff, M., 2016. Objective assessment of the contribution of the RECOPECA network  
850 to the monitoring of 3D coastal ocean variables in the Bay of Biscay and the English  
851 Channel. *Ocean Dynamics*, 66(4), 567-588, <http://dx.doi.org/10.1007/s10236-016-0938-y>.
- 852 Le Hénaff, M., P. De Mey and Marsaleix, P., 2009. Assessment of observational networks with  
853 the Representer Matrix Spectra method – Application to a 3-D coastal model of the Bay of  
854 Biscay. *Ocean Dynamics*, 59, 3-20, DOI 10.1007/s10236-008-0144-7.
- 855 Leonardo, L., R. Matear, M. Soja-Woźniak, D. J. Suggett, D. J. Hughes, M. E. Baird, M. A.  
856 Doblin, 2021. Modelling the impact of phytoplankton cell size and abundance on inherent  
857 optical properties (IOPs) and a remotely sensed chlorophyll-a product, *Journal of Marine*  
858 *Systems*, 213, 103460, <https://doi.org/10.1016/j.jmarsys.2020.103460>.
- 859 Madec, G., 2012. Nemo ocean engine, Tech. rep., NEMO team.
- 860 Mattern, J.P., Edwards, C.A., and Moore, A.M., 2018. Improving Variational Data  
861 Assimilation through Background and Observation Error Adjustments. *Mon. Wea. Rev.*,  
862 146, 485–501, <https://doi.org/10.1175/MWR-D-17-0263.1>.
- 863 Oke, P. R., and P. Sakov, 2008. Representation Error of Oceanic Observations for Data  
864 Assimilation. *J. Atmos. Oceanic Technol.*, 25, 1004–1017,  
865 <https://doi.org/10.1175/2007JTECHO558.1>.
- 866 Pujol, M.-I., Faugère, Y., Taburet, G., Dupuy, S., Pelloquin, C., Ablain, M., and Picot, N.,  
867 2016. DUACS DT2014: the new multi-mission altimeter data set reprocessed over 20 years,  
868 *Ocean Sci.*, 12, 1067-1090, doi:10.5194/os-12-1067-2016.
- 869 Quattrocchi, G., P. De Mey, N. Ayoub, V. Vervatis, C.-E. Testut, G. Reffray, J. Chanut and Y.  
870 Drillet, 2014. Characterisation of errors of a regional model of the Bay of Biscay in response

- 871 to wind uncertainties: a first step toward a data assimilation system suitable for coastal sea  
872 domains. *Journal of Operational Oceanography*, Volume 7, Number 2, August 2014, pp. 25-  
873 34(10).
- 874 Sakov, P., G. Evensen and L. Bertino, 2010. Asynchronous data assimilation with the EnKF.  
875 *Tellus A.* 62(1): 24-29, <http://dx.doi.org/10.1111/j.1600-0870.2009.00417.x>.
- 876 Sieburth, John McN. Smetacek, Victor Lenz, Jürgen, 1978. Pelagic ecosystem structure:  
877 Heterotrophic compartments of the plankton and their relationship to plankton size  
878 fractions, *Limnology and Oceanography*, 23, doi: 10.4319/lo.1978.23.6.1256.
- 879 Simon, E., and Bertino, L., 2009. Application of the Gaussian anamorphosis to assimilation in  
880 a 3-D coupled physical-ecosystem model of the North Atlantic with the EnKF: a twin  
881 experiment. *Ocean Sci.* 5, 495-510.
- 882 Skakala, J., Ford, D., Brewin, R. J. W., Mc Ewan, R., Kay, S., Taylor, B., de Mora, L., and  
883 Ciavatta, S., 2018. The assimilation of phytoplankton functional types for operational  
884 forecasting in the North-West European Shelf. *Journal of Geophysical research - Oceans*,  
885 123 (8), 5230-5247, doi: 10.1029/2018JC014153.
- 886 Song, H., Edwards, C.A., Moore, A.M., Fiechter, J., 2016. Data assimilation in a coupled  
887 physical-biogeochemical model of the California current system using an incremental  
888 lognormal 4-dimensional variational approach: Part 3—Assimilation in a realistic context  
889 using satellite and in situ observations. *Ocean Model.* 106, 159–172.  
890 <https://doi.org/10.1016/j.ocemod.2016.06.005>.
- 891 Vervatis, V. D., C.E. Testut, P. De Mey, N. Ayoub, J. Chanut, and G. Quattrocchi, 2016. Data  
892 assimilative twin-experiment in a high-resolution Bay of Biscay configuration: 4D EnOI  
893 based on stochastic modelling of the wind forcing. *Ocean Modelling*, 100, 1-19,  
894 <http://dx.doi.org/10.1016/j.ocemod.2016.01.003>.
- 895 Vervatis, V. D., De Mey-Frémaux, P., Ayoub, N., Karagiorgos, J., Ghantous, M., Kailas, M.,  
896 Testut, C.-E., and Sofianos, S., 2021. Assessment of a regional physical-biogeochemical  
897 stochastic ocean model. Part 1: ensemble generation, *Ocean modell.*, under review.

898 **Table 1.** Stochastic model ensembles (cf. companion article Part 1 by *Vervatis et al.*, 2021).

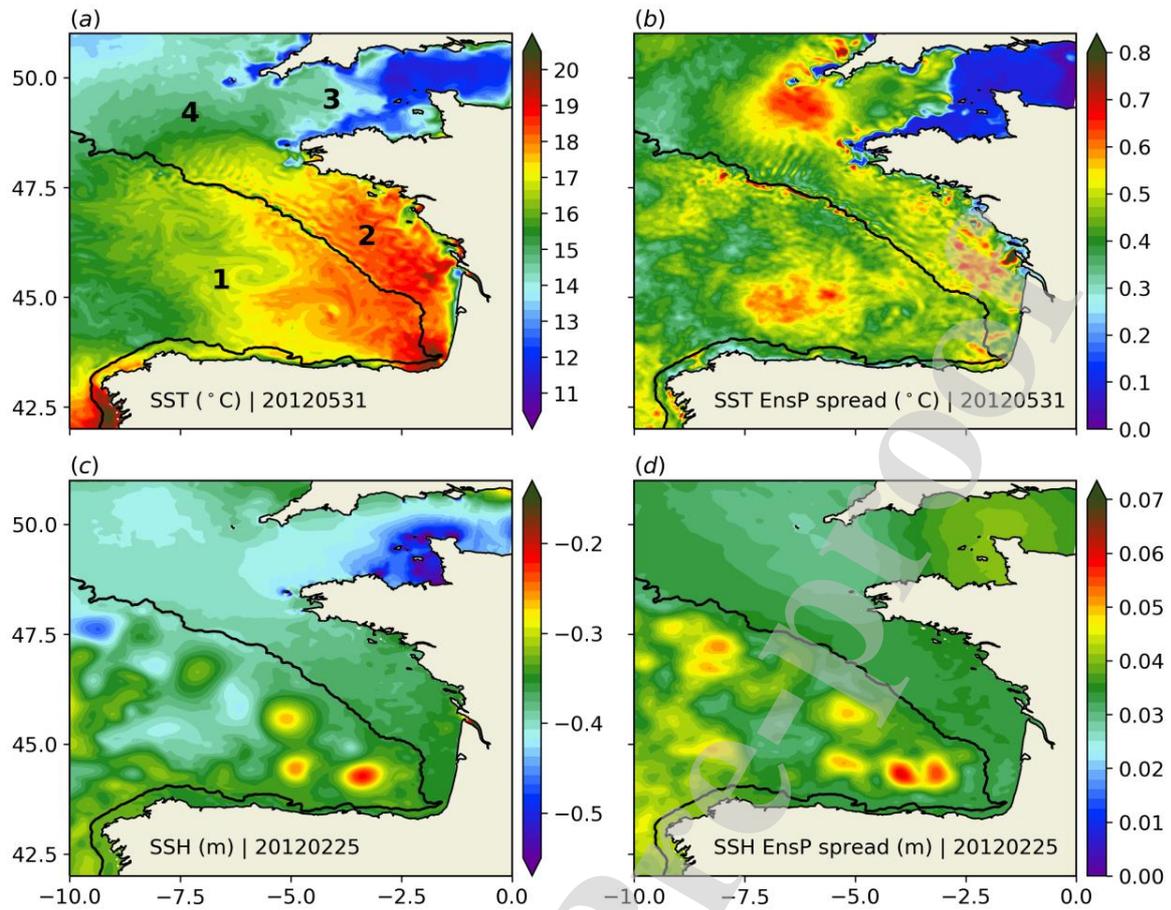
Ensemble	SPPT/SPP-AR(1)		
EnsP	Physics	Atm. forcing Parameters	$U_{air}, T_{air}, SLP$ $c_d, c_e, c_h, c_b$
EnsB	Biogeochemistry	sources-minus-sinks of 24 prognostic variables $SMS(C)$	
EnsPB	Physics & Biogeochemistry	EnsP & EnsB	

899 abbreviations: SPPT - stochastic perturbed parameterized tendencies; SPP - stochastic  
900 perturbed parameters; AR(1) - first-order autoregressive processes;  $U_{air}$  - wind velocities;  $T_{air}$   
901 - air temperature;  $SLP$  - sea level pressure;  $c_d, c_e, c_h$  - wind drag and turbulent coefficients;  $c_b$   
902 - bottom drag;  $SMS(C)$  - sources minus sinks of biogeochemical tracers  $C$ .

903 **Table 2.** Observational networks.

CMEMS Product Identifiers ( <a href="http://marine.copernicus.eu/">http://marine.copernicus.eu/</a> )*		Error
OSTIA SST L4 gridded 0.05°	SST_GLO_SST_L4_NRT_OBS_010_001	0.5 °C
NWS SST L4 gridded 0.04°	SST_ATL_SST_L4_REP_OBS_010_026	
SLA L3 along-track 14 km	SEALEVEL_GLO_PHY_L3_REP_OBS_008_062	0.05 m
SLA L4 gridded 0.25°	SEALEVEL_GLO_PHY_L4_REP_OBS_008_047	
Chlorophyll <i>a</i> gridded 0.05°	Phytoplankton Functional Types - PFT ( <i>Brewin et al.</i> , 2017)	0.3 mg/m <sup>3</sup>

904 \*All datasets are provided at daily frequency.



905

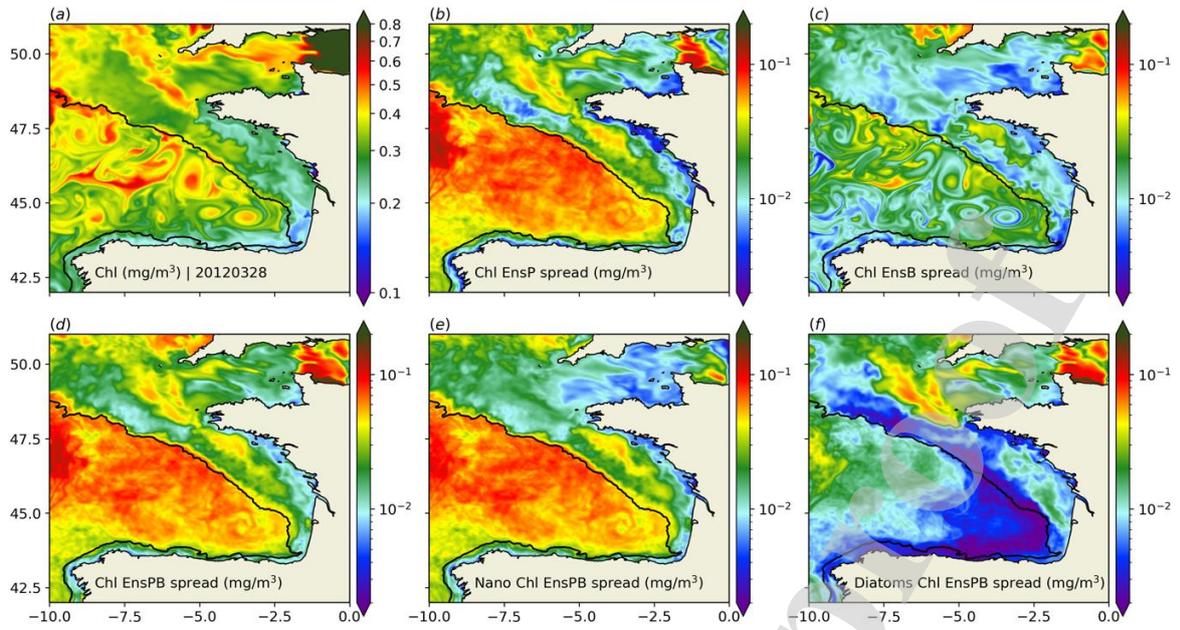
906

907

908

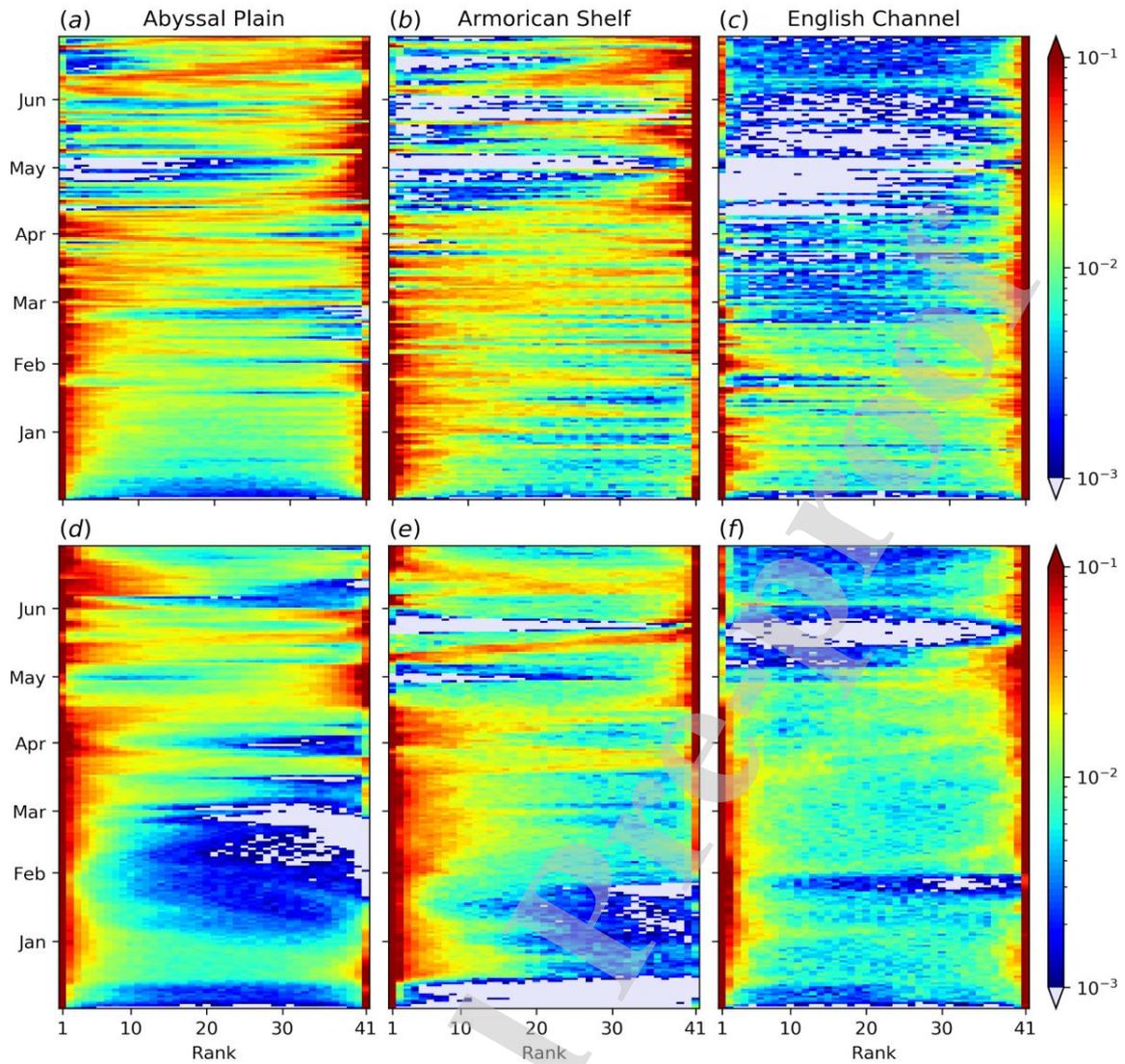
909

**Figure 1** (a-b) Deterministic run SST and model ensemble EnsP spread ( $^{\circ}\text{C}$ ) on May 31, 2012. (c-d) Same for SSH ( $m$ ) on February 25, 2012. Fig. 1a depicts the characteristic areas in the Bay of Biscay discussed in the text: 1-abyssal plain, 2-Armorican shelf, 3-English Channel, 4-Celtic shelf. Black line denotes the 200 m isobaths.



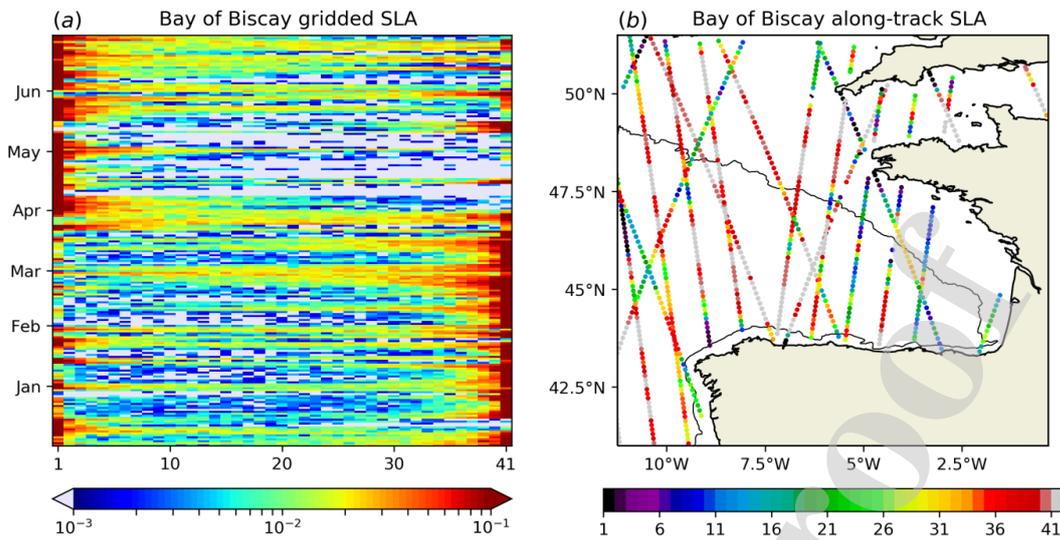
910

911 **Figure 2** (a-d) Total chlorophyll deterministic run and model ensembles spreads EnSP,  
 912 EnSB, EnSPB ( $mg/m^3$ ) on March 28, 2012. (e-f) Same with (d) model ensemble spread EnSPB  
 913 for nanophytoplankton and diatoms chlorophyll ( $mg/m^3$ ) respectively.



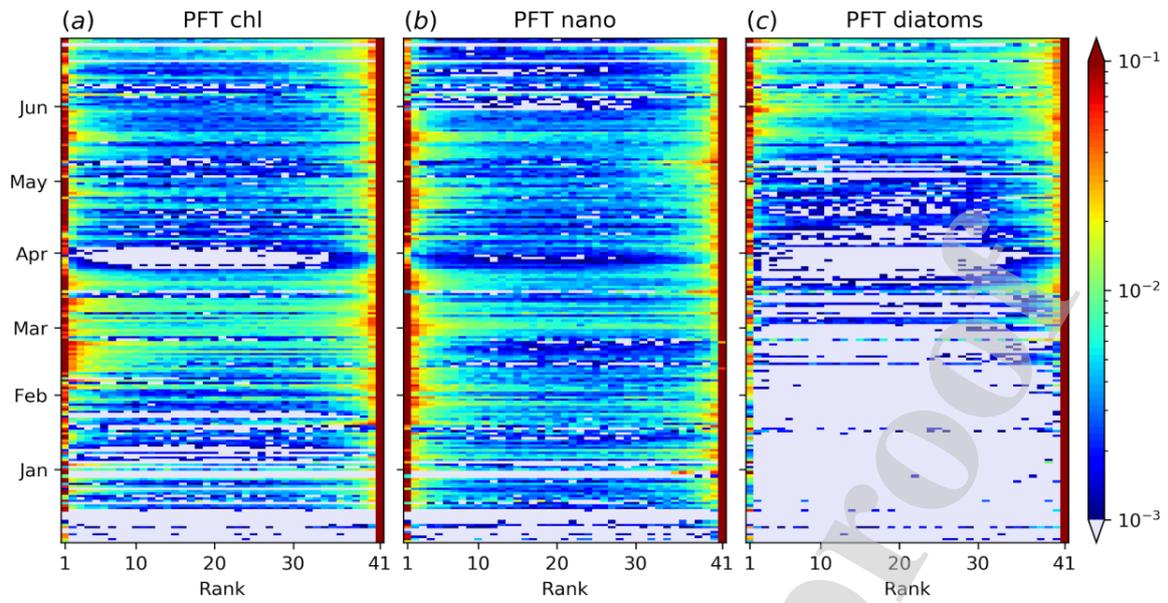
914

915 **Figure 3** Hovmöller plot of rank histograms (x-axis: rank, y-axis: time, colorbar: *pdf*)  
 916 between ensemble EnsP and SST L4 datasets: (a-c) OSTIA, (d-f) NWS. Regional consistency  
 917 is checked for the (left to right): abyssal plain, Armorican shelf and English Channel. Rank  
 918 histogram flatness is verified for *pdf* values at  $\frac{1}{41} \sim 0.024$  (green-yellow in colorbar).



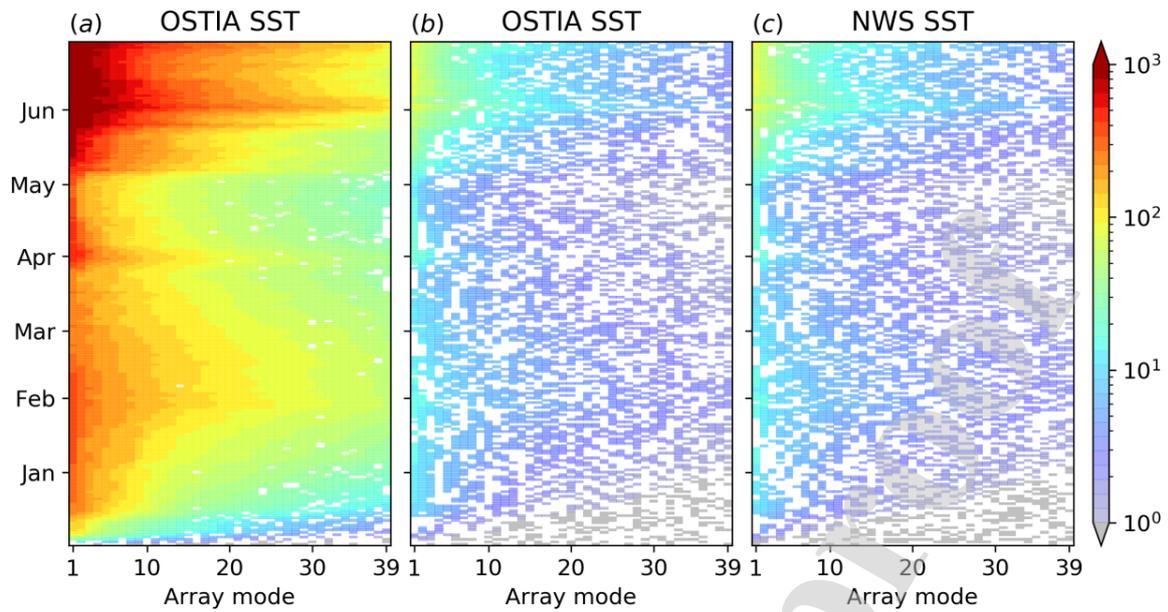
919

920 **Figure 4** (a) Hovmöller plot of rank histograms (x-axis: rank, y-axis: time, colorbar: *pdf*)  
 921 between ensemble EnsP and SLA L4 gridded dataset for the Bay of Biscay, and (b) map of  
 922 ranks for along-track SLA L3 dataset with respect to the ensemble EnsP (colorbar: rank),  
 923 starting on February 25, 2012, and for 10 consecutive days, including ascending and  
 924 descending tracks in the domain.



925

926 **Figure 5** Hovmöller plot of rank histograms (x-axis: rank, y-axis: time, colorbar: *pdf*)  
927 between ensemble EnsPB and PFT dataset: (a) total chlorophyll, (b) nanophytoplankton and  
928 (c) diatoms chlorophyll.



929

930

931

932

933

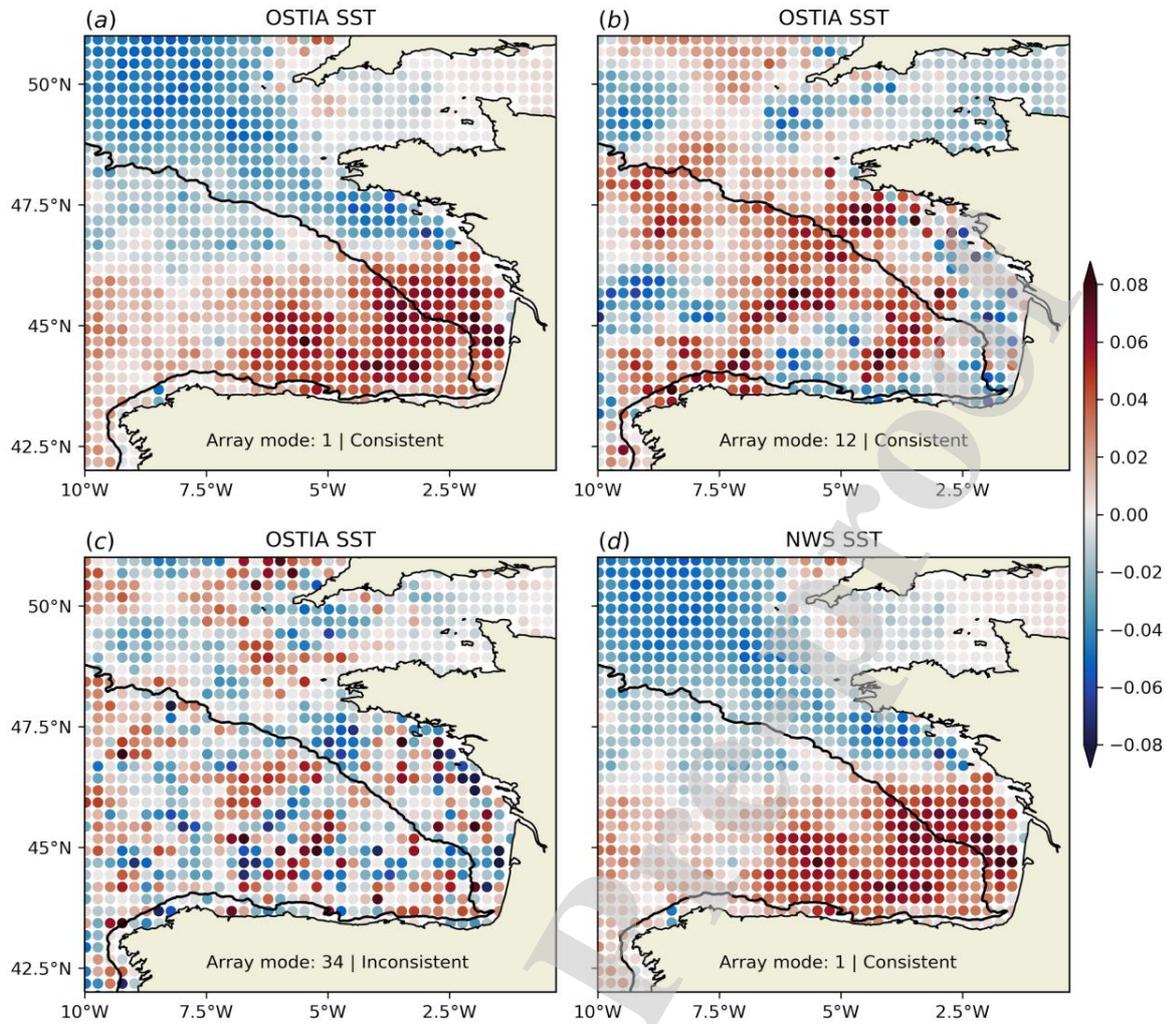
934

935

936

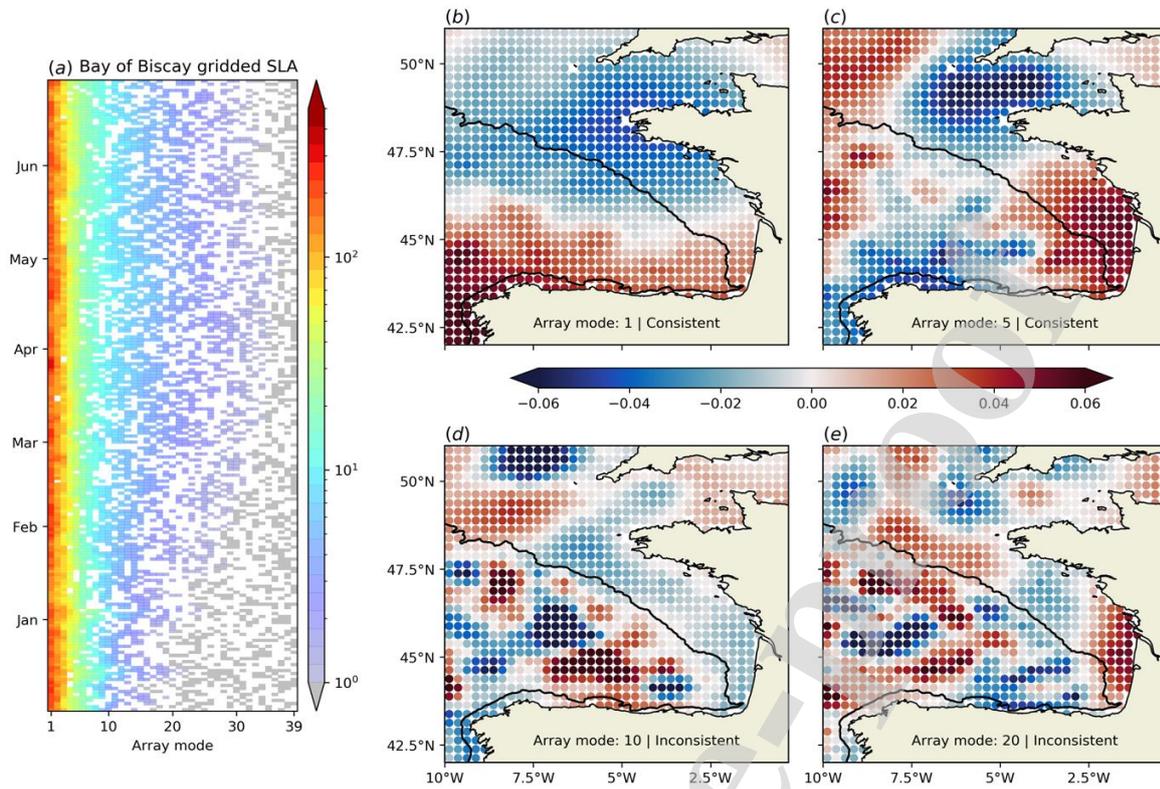
937

**Figure 6** Hovmöller plot of variations in time of SST array mode spectra vs. modal rank (i.e. the rank of the array mode), including an EnsP consistency check against (a-b) OSTIA and (c) NWS SST data. Colorbar: array mode spectra as of criterion ArM1; eigenvalues smaller than one (the observational noise floor in array space; grey pixels) denote error modes marginally detectable by the array. White pixels depict inconsistent array modes as per criterion ArMCA1 (the higher modes are also mostly inconsistent as per criterion ArM1). (a) No data subsampling, (b) one-fifth subsampling rate (i.e. one data point every fifth OSTIA point retained), and (c) one-sixth subsampling rate.



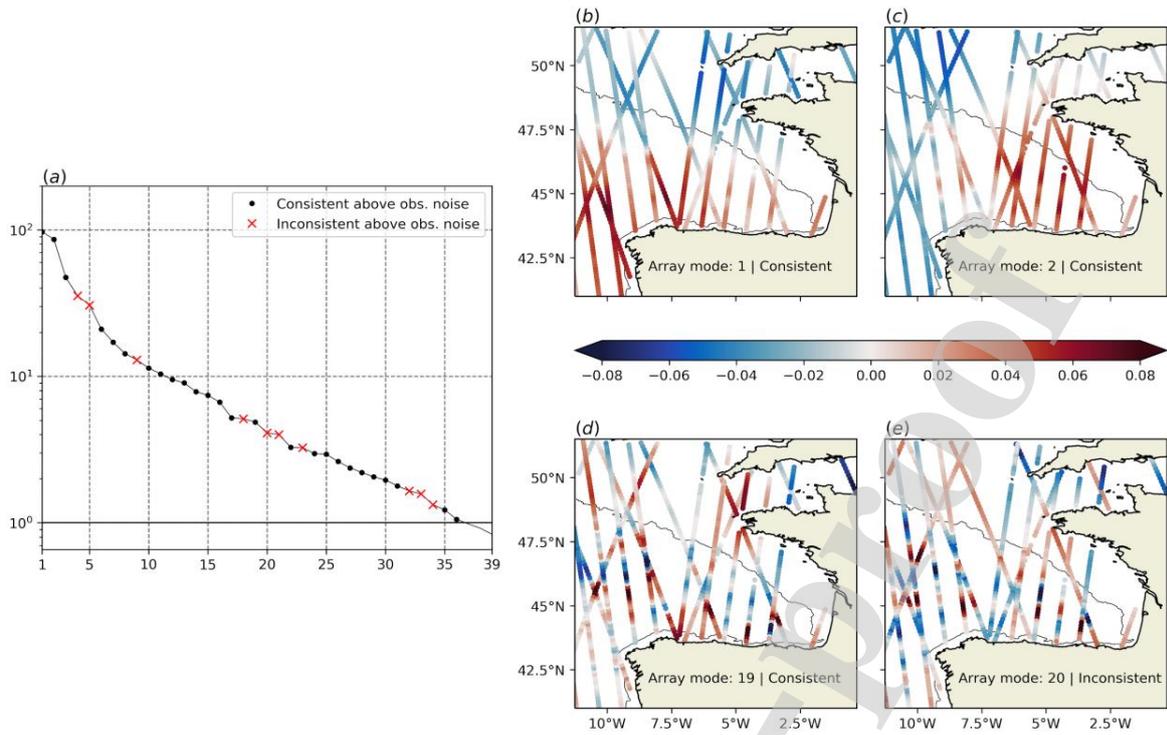
938

939 **Figure 7** (a-c) Array modes 1, 12 and 34 for EnsP as “seen” from OSTIA SST on May  
 940 31, 2012, using one-fifth subsampling rate, and (below as text) corresponding results of the  
 941 ArMCA1 consistency criterion using OSTIA SST observations; the 1st and 12th array modes  
 942 are consistent; array mode 34 is inconsistent. (d) Same as (a) for NWS SST and one-sixth  
 943 subsampling rate. Colorbar: array mode amplitude (no units) on (subsamped) data grid.



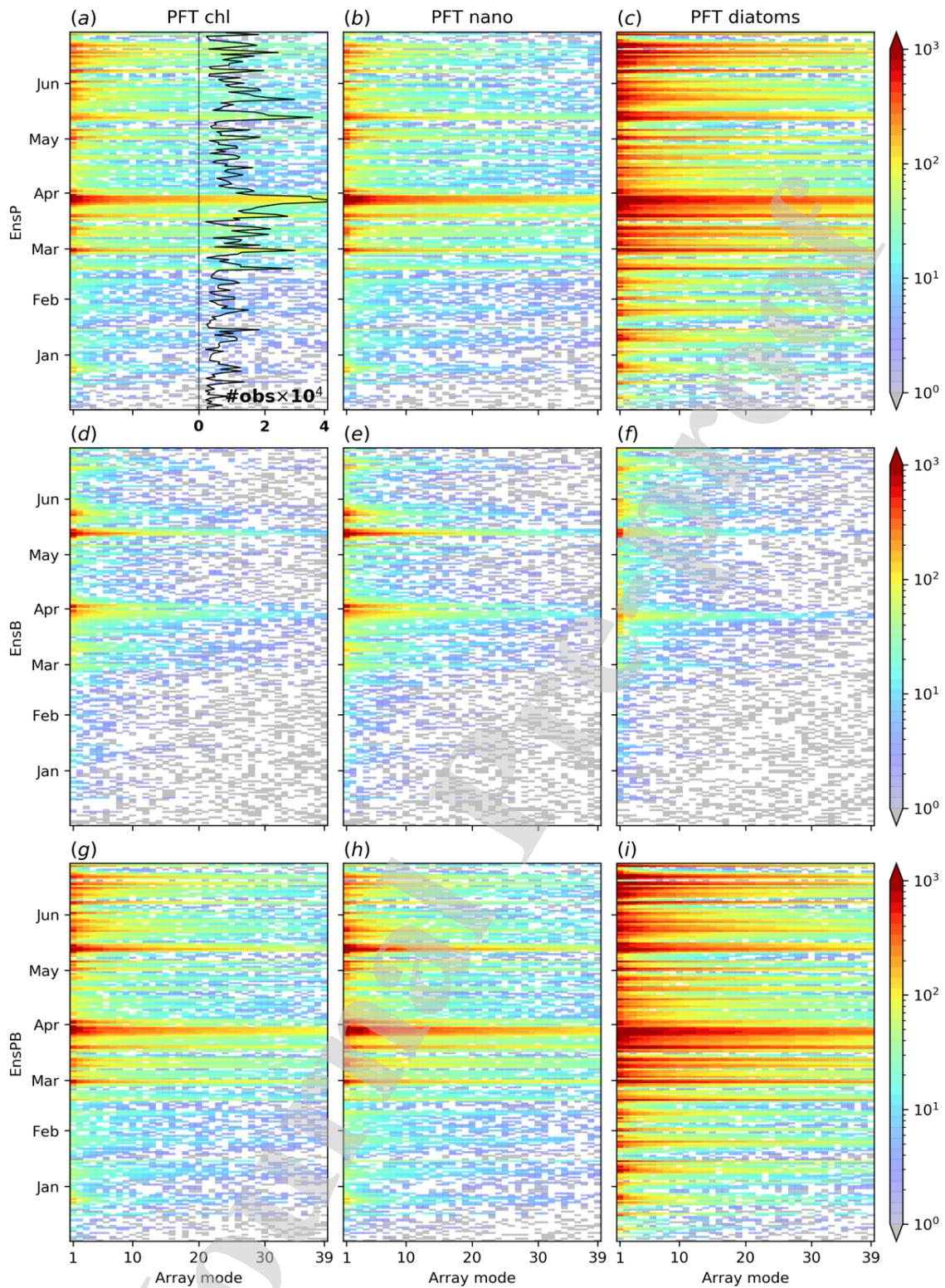
944

945 **Figure 8** (a) Hovmöller plot of variations in time of gridded SLA array mode spectra vs.  
 946 modal rank, including an EnsP consistency check against gridded SLA data. (b-e) Examples of  
 947 consistent and inconsistent gridded SLA array modes on April 30, 2012 for EnsP. Colorbars  
 948 and units as in Figs. 6 and 7.



949

950 **Figure 9** (a) Along-track SLA L3 array mode spectrum for the same period as in Fig. 4b,  
 951 considering all tracks that would be assimilated in a 10-day assimilation cycle. Black and red  
 952 markers for eigenvalues above the observational noise floor (= 1 in array space) denote  
 953 consistent and inconsistent modes respectively. (b-c) Multi-track network consistent array  
 954 modes of ranks 1 and 2 respectively (no units; all times “flattened”, i.e. 2D representation of  
 955 array modes including time space); (d-e) same as (b-c) for a higher consistent array mode of  
 956 rank 19 and for an inconsistent array mode of rank 20 respectively. Colorbar and units as in  
 957 Figs. 6 and 7.



958

959

960

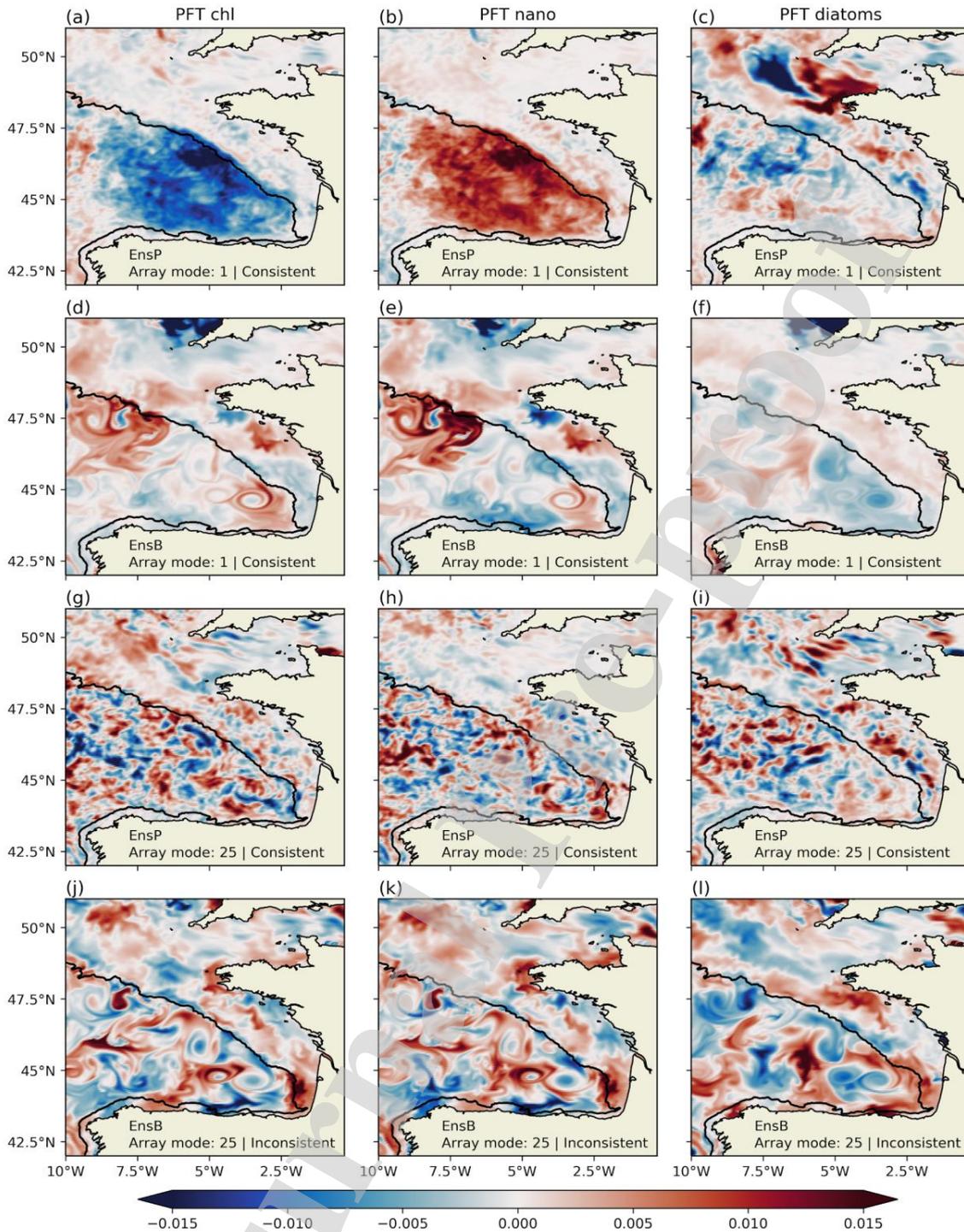
961

962

963

964

**Figure 10** (a-c) Hovmöller plots of variations in time of PFT array mode spectra vs. modal rank, including EnsP consistency checks against PFT data for total chlorophyll, nanophytoplankton, and diatoms respectively as shown; (d-f) same as (a-c) for EnsB; (g-i) same as (a-c) for EnsPB. Chlorophyll distributions have been log transformed. Colorbars and units as in Fig. 6. Fig. 10a superimposed black line: number of data with largest values corresponding to full data coverage of the domain i.e. #obs  $\sim (4 \cdot 10^4)$ .



965

966 **Figure 11** Examples of array modes for EnsP and EnsB on March 28, 2012, as “seen” from  
 967 PFT data points, and (below as text) corresponding results of the ArMCA1 consistency  
 968 criterion using PFT data: (a-c) array modes at rank 1 using EnsP (all consistent); (d-f)  
 969 array modes at rank 1 using EnsB (all consistent, but different); (g-i) array modes at rank 25  
 970 using EnsP (all consistent); (j-l) array modes at rank 25 using EnsB (all inconsistent). PFT gridded  
 971 data (gappy in general) have no gaps on that particular date in the Bay of Biscay. (left column)  
 972 Total chlorophyll; (center) nanophytoplankton; (right column) diatoms. No data subsampling.  
 973 Array modes calculated in log array space. Colorbar and units as in Fig. 7.

**Highlights**

- Model-data biases and model underspread revealed in rank histograms
- Consistent SST array modes at large scales and at small-scale
- Along-track array modes showed useful consistent information at the mesoscale
- Consistent PFT array modes at small-scale perturbing physics
- Additional error processes active in the model for inconsistent configurations

**Authors' contribution statement using CRediT**

- 1) Vassilios D. Vervatis: conceptualization, methodology, software, visualization, formal analysis, investigation, writing – original draft, writing – review and editing.
- 2) Pierre De Mey-Frémaux: conceptualization, methodology, software, formal analysis, writing – original draft, writing – review and editing.
- 3) Nadia Ayoub: conceptualization, formal analysis, writing – original draft, writing – review and editing.
- 4) John Karagiorgos: software, visualization, formal analysis.
- 5) Stefano Ciavatta: data curation, writing – review and editing.
- 6) Robert J. W. Brewin: data curation, writing – review and editing.
- 7) Sarantis Sofianos: conceptualization, resources.

The corresponding author,

Vassilios Vervatis

**Declaration of interests**

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

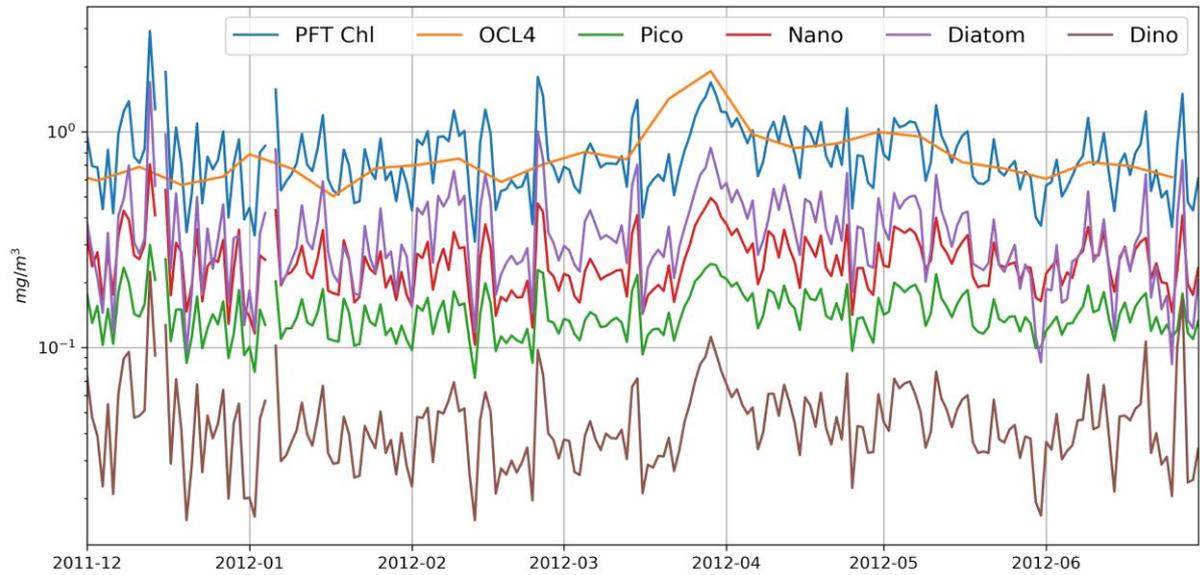
Journal Pre-proof

## 974 **Supplementary material**

975 Fig. S1 presents the temporal variability of the PFT data in the Bay of Biscay, during the period  
976 December, 2011 to June, 2012. We also show the ocean colour L4 total chlorophyll product  
977 used in the Part 1 article (*Vervatis et al.*, 2021) and we verify that concentrations are of the  
978 same order to those of the PFT total chlorophyll. According to this, we categorised the four  
979 satellite PFT (pico, nano, diatoms and dino) into the two broad size groups in PISCES (nano  
980 or diatoms), in a manner most representative, ensuring that the total biomass (chlorophyll) from  
981 the model and satellite data can be compared like-for-like.

982 PFT diatoms and nanophytoplankton contribute together approximately more than 80% in total  
983 chlorophyll, whereas picoplankton contributes at about 10% and dinoflagellates less than 10%  
984 (Fig. S1). PFT diatom chlorophyll concentration is an order of magnitude larger compared with  
985 dinoflagellates and nanophytoplankton is about three times larger than picoplankton  
986 chlorophyll concentration (Fig. S1). In Fig. S2, we show the spatial distribution of the four  
987 satellite PFT and the total chlorophyll during the peak of the spring bloom on March 28, 2012.  
988 We confirm that the satellite micro class (i.e. diatoms and dino) is driven primarily by diatoms,  
989 far more abundant in the satellite data, with the two functional types being highly correlated in  
990 spatial. Fig. S3 presents scatter plots of combined vs. non-combined PFT chlorophyll, verifying  
991 the close relationship between functional types in a size class-based approach.

992 We also present results from one-on-one comparisons between model and PFT data, as opposed  
993 to the size class-based categorization merging different functional types. Figure S4 shows  
994 Hovmöller plot of rank histograms between EnsPB and PFT, in the same way as Figs. 5b-c,  
995 with one main difference: in Figs. S4a-b we do not combine the nano functional type with pico,  
996 nor we combine diatoms with dino. Rank histogram results for the nano class are degraded  
997 when pico and nano PFT data are not combined together in late-winter and early-spring when  
998 a primary bloom occurs (Fig. 5b vs. Fig. S4a). The latter may suggest that PISCES nano can  
999 be representative of a broader phytoplankton community, accounting also for smaller size  
1000 classes. Rank histogram results are almost identical throughout the whole period for the micro  
1001 class, regardless of whether dino and diatoms are combined together or not (Fig. 5c vs. Fig.  
1002 S4b). Model-data one-on-one array mode consistency results (not shown) are in practice  
1003 indistinguishable by visual inspection with the results presented in Figs. 10 and 11, confirming  
1004 the validity of the size class-based approach.

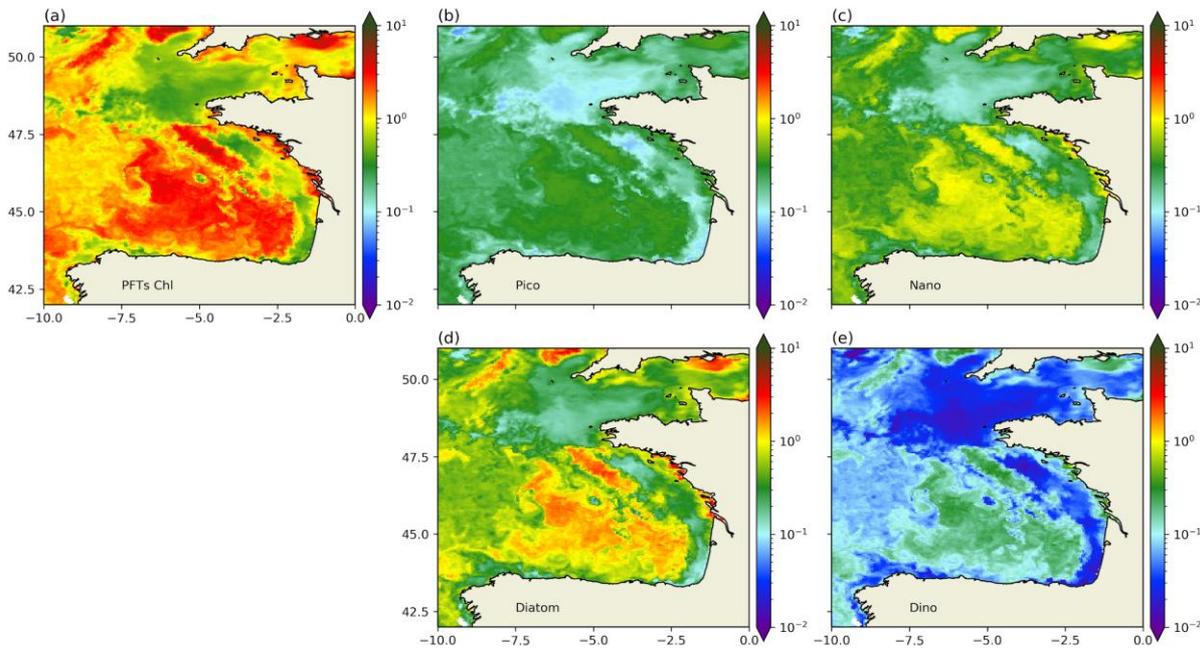


1005

1006

1007

**Figure S1** Ocean colour L4 (8-day frequency) and PFT (daily) chlorophyll concentration ( $mg/m^3$ ) in the Bay of Biscay from December, 2011 to June, 2012.

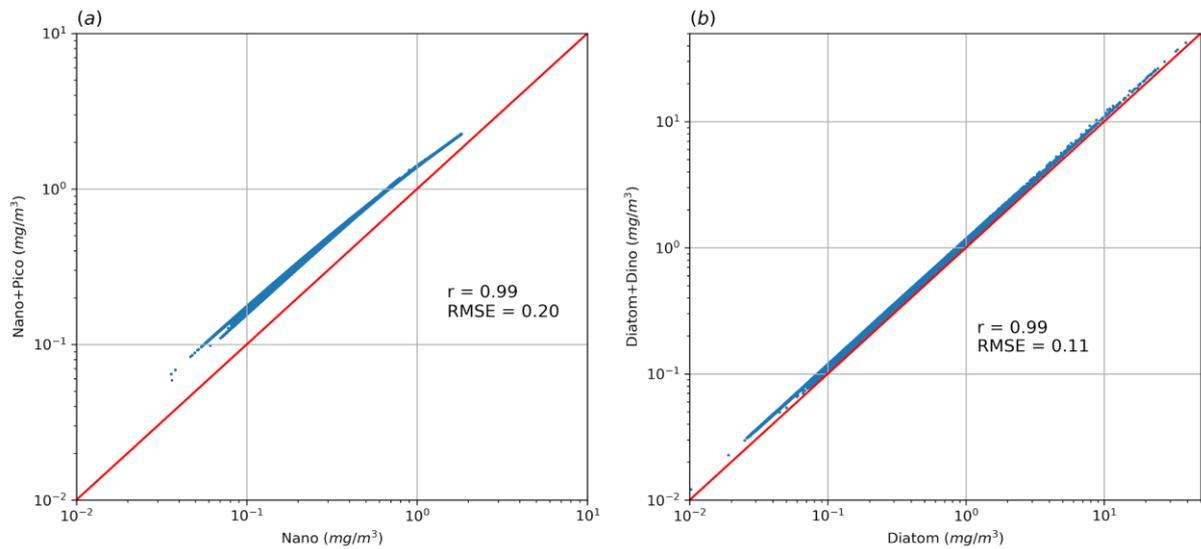


1008

1009

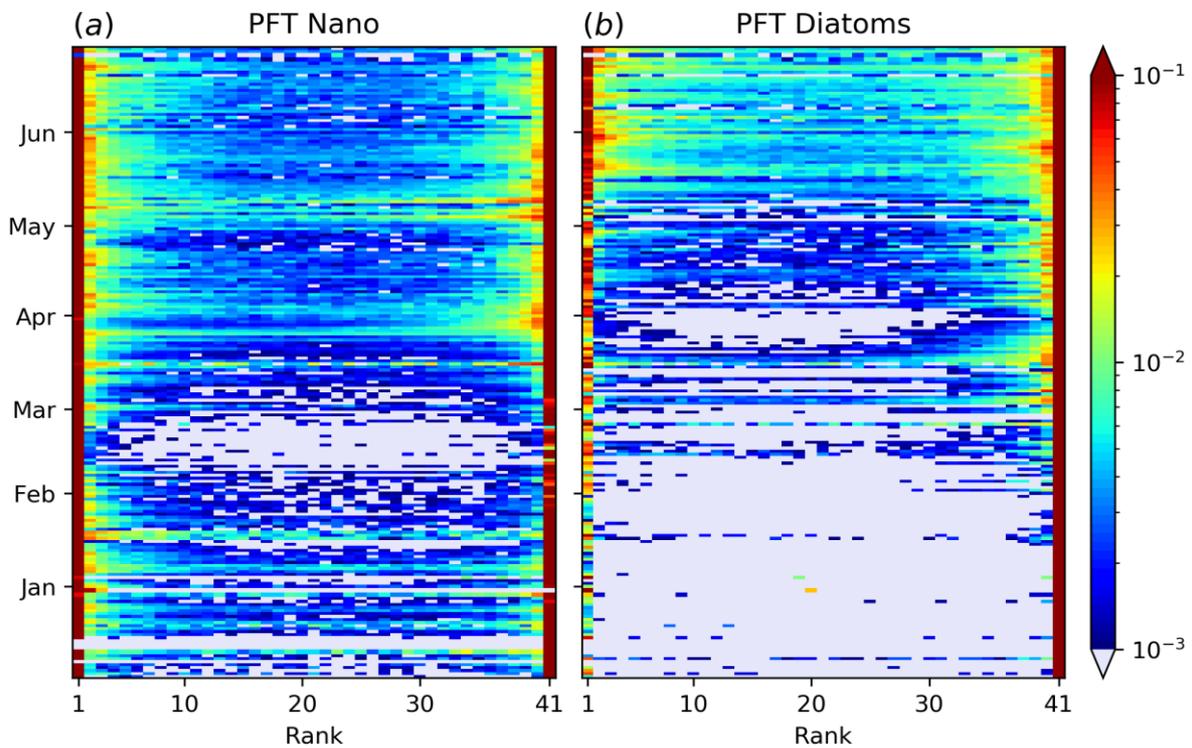
1010

**Figure S2** Spatial distribution of (a) PFT total chlorophyll and (b) pico, (c) nano, (d) diatoms, (e) dino concentrations ( $mg/m^3$ ) on March 28, 2012.



1011

1012 **Figure S3** Scatter plots of chlorophyll concentrations in  $mg/m^3$  on March 28, 2012: (a)  
 1013 PFT (nano and pico) vs. only nano, (b) PFT (diatoms and dino) vs. only diatoms.  $r$  is the  
 1014 correlation coefficient, RMSE the root mean square error ( $mg/m^3$ ) and with red the 1:1 line.



1015

1016 **Figure S4** Hovmöller plot of rank histograms (same as in Figs. 5b-c) between EnsPB and  
 1017 PFT (a) nano not combined with pico, and (b) diatoms not combined with dino.