

# Filtering remotely sensed chlorophyll concentrations in the Red Sea using a space-time covariance model and a Kalman Filter

Denis Dreano<sup>a</sup>, Bani Mallick<sup>b</sup>, Ibrahim Hoteit<sup>a,\*</sup>

<sup>a</sup>*Computer, Electrical and Mathematical Sciences and Engineering Division, King Abdullah University of Science and Technology*

<sup>b</sup>*Department of Statistics, Texas A&M University*

---

## Abstract

A statistical model is proposed to filter satellite-derived chlorophyll concentration from the Red Sea, and to predict future chlorophyll concentrations. The seasonal trend is first estimated after filling missing chlorophyll data using an Empirical Orthogonal Function (EOF)-based algorithm (Data Interpolation EOF). The anomalies are then modeled as a stationary Gaussian process. A method proposed by Gneiting [1] is used to construct positive-definite space-time covariance models for this process. After choosing an appropriate statistical model and identifying its parameters, Kriging is applied in the space-time domain to make a one step-ahead prediction of the anomalies. The latter serves as the prediction model of a reduced-order Kalman filter, which is applied to assimilate and predict future observations of chlorophyll concentrations. The proposed method decreases the root mean square

---

\*Corresponding author

Email address: [ibrahim.hoteit@kaust.edu.sa](mailto:ibrahim.hoteit@kaust.edu.sa) (Ibrahim Hoteit)

(RMS) prediction error by about 11% compared with the seasonal prediction.

---

## **1. Introduction**

The Red Sea is an elongated basin situated between the Asian and the African shelves, connected to the Mediterranean Sea in the north through the Suez Canal, and to the Gulf of Aden in the south through the Strait of Bab el Mandeb [2]. It is one of the warmest and most saline seas in the world, with a rich ecosystem that has adapted to these extreme conditions [3]. This unique natural resource is however threatened by an abrupt increase of temperature since 1994 [4].

Phytoplankton are small, unicellular, photosynthetic algae. They are the primary producers for marine ecosystems, and at the base of the marine food chain. Phytoplankton concentration is therefore important for fisheries [5]. By fixing atmospheric CO<sub>2</sub> and sinking to form sediment at the bottom of the sea, phytoplankton also acts as a biological pump. This phenomenon is crucial to understanding climate change [6]. Thus modeling and predicting changes in chlorophyll concentrations have many purposes. In the Red Sea, phytoplankton is particularly important to the extensive coral reefs along its shores.

The Red Sea is generally deficient in major inorganic nutrients, and its productivity is relatively low [7, 8]. The high productivity observed in the southern Red Sea is attributed to the intrusion of nutrient-rich waters from the Gulf of Aden [7, 8]. Red Sea chlorophyll concentrations follow seasonal

22 patterns, with a winter bloom following a weak summer productivity. Con-  
23 siderable interannual variability in chlorophyll concentrations has also been  
24 observed [9]. However, the Red Sea ecosystem has not yet been fully ex-  
25 plored, and very few in- situ measurements have been conducted in its basin  
26 [10], increasing the need for remotely sensed data. Satellite observations of  
27 chlorophyll concentrations have been shown to be reliable datasets to study  
28 the primary productivity of the oceans [11] and they constitute the basis of  
29 several studies in the Red Sea [7, 10, 9].

30 Modeling and predicting changes in phytoplankton concentrations is chal-  
31 lenging. It requires the coupling of an ecological model with a hydrological  
32 model, which provides the physical forcing that influences the ecology. A  
33 broad range of models has been developed by the marine ecosystem research  
34 community, from the simple NPZ model with only nutrients, phytoplankton  
35 and zooplankton as state variables, to much more complex models. NPZ  
36 models have been implemented with various degrees of success on very dif-  
37 ferent regions [12]. However, there are cases when distinguishing between  
38 specific groups of phytoplankton can be useful. Examples are the study  
39 of export, sinking or climate feedback [12]. The European Regional Seas  
40 Ecosystem Model (ERSEM) [13] is an example of a much more sophisticated  
41 ecosystem model, initially designed for simulations of the North Sea. It dis-  
42 tinguishes functional phytoplankton and zooplankton groups, and models the  
43 complete cycling of different nutrient groups and O<sub>2</sub> and CO<sub>2</sub>, and includes  
44 the effect of higher trophic groups. It has recently been implemented in the

45 Red Sea [8]. Configuring ERSEM and coupling it with a physical ocean  
46 model is, however, delicate and requires considerable efforts and expertise  
47 [14]. A more important issue with this class of models is the number of  
48 empirical equations that governs its dynamics and the number of involved  
49 ecological variables and parameters (over 50 for ERSEM [15]). This makes  
50 such models difficult to calibrate and validate, since there are usually not  
51 enough observations to constrain the large number of parameters [12].

52 An alternative approach is to follow a statistical framework to model the  
53 space-time evolution of chlorophyll concentrations. Statistical methods have  
54 not been widely used in this field, which has previously relied on time-series  
55 observations. Artificial neural networks were applied to forecasting algal  
56 blooms in freshwater and marine systems [16, 17], and generalized additive  
57 models have been use in finding explanatory variables for the chlorophyll  
58 concentrations in the Pagasitikos Gulf and the subartic North Atlantic [18,  
59 19].

60 Geostatistical spatio-temporal models are extensions of the spatial clas-  
61 sical geostatistical methods [20]. These methods consider space-time data as  
62 the realization of a Gaussian process, from which a mean and a covariance  
63 function can be estimated. Although, in most applications, such a stochas-  
64 tic modeling approach has no rigorous scientific basis, geostatistical methods  
65 may capture some patterns in the data and avoid the difficulties of dynam-  
66 ical models [21, 20]. These methods are widely employed in meteorology to  
67 model the surface temperature over land and oceans [22, 23], in an ecological

68 context to study moth populations [24], and to characterize soil and pollution  
69 [25, 26].

70 The first proposed spatio-temporal geostatistical model was based on sep-  
71 arable covariance functions [20]. Such functions can either be the product, or  
72 the sum, of a purely temporal covariance model and a purely spatial covari-  
73 ance model. These covariance models are convenient but have non-physical  
74 properties, limiting their use in many situations [20]. As a result, significant  
75 research has been recently conducted to construct nonseparable families of  
76 covariance functions. In [27, 1], the authors proposed methods to construct  
77 families from known temporal and spatial covariance functions. The method  
78 proposed in [1] is appealing because of its modularity and interpretability.  
79 We adopte it here to model the anomaly fields of chlorophyll concentration  
80 in the Red Sea.

81 Once the mean and covariance functions have been estimated from avail-  
82 able satellite data after missing data have been filled with an Empirical Or-  
83 thogonal Function (EOF)-based method (DINEOF [28, 29]), it is possible to  
84 condition on present observations to improve the forecast of future chloro-  
85 phyll concentration measurements via Kriging [20]. However, this requires  
86 the inversion of the spatial observation autocovariance matrix with a 0 tem-  
87 poral. In the context of remotely sensed two-dimensional (2D) fields, the  
88 number of observations may be close to the number of variables. For a  
89 large area and/or high resolution, this inversion can be computationally pro-  
90 hibitive. Moreover, in the case of large regions of missing observations, as

91 in the southern Red Sea, the prediction will be biased at these locations  
92 [30]. This problem is tackled here by identifying the Kriging operator as  
93 an evolution matrix in a state-space context and then using a reduced-order  
94 Kalman filter [31] for the filtering and prediction of the observations. This  
95 method is computationally efficient; it implicitly uses past observations in  
96 the estimation process; and it provides a prediction of the entire state at  
97 any point in time. The reduced-order Kalman filter significantly reduces the  
98 computational burden of the traditional Kalman filter by projecting the filter  
99 covariance matrices on a fixed low-rank basis. Since the anomalies are mod-  
100 eled with a space-time Gaussian process, they are correlated in time. The  
101 formulation of the reduced-order Kalman filter is therefore expanded to take  
102 into account a colored model noise [32].

103 The paper is organized as follows. Section 2 describes the satellite data.  
104 Section 3 discusses the methodology for the data filling and geostatistical  
105 modeling. Section 4 derives the state-space formulation and introduces the  
106 space-time Kalman filtering problem and its solution. The experimental  
107 setup and numerical results are presented and discussed in section 4. Con-  
108 cluding remarks are provided in section 5.

## 109 **2. Data and preprocessing**

110 Satellite data provide chlorophyll (CHL) concentrations with a spatial and  
111 temporal resolution not achievable with in situ observations, making them  
112 particularly relevant to the Red Sea, where very few in situ data collection

113 are conducted.

114 Level-3 mapped data from the NASA SeaWiFS (Sea-Viewing Wide Field-  
115 of-View Sensor) satellite sensor are used in this study. The dataset is publicly  
116 available at <http://oceancolor.gsfc.nasa.gov>. In this study, we use the  
117 9km resolution mapped weekly averages from January 1998 to December  
118 2007 (460 time steps). At each time step, a  $133 \times 188$  pixel map is avail-  
119 able for a domain extending from longitudes between  $33^{\circ}\text{E}$  and  $44^{\circ}\text{E}$  and  
120 latitudes between  $12^{\circ}\text{N}$  and  $28^{\circ}\text{N}$ , of which 5635 pixels correspond to actual  
121 Red Sea surface (see Figure 1(a)). A log-transformation was applied in or-  
122 der to obtain an approximately Gaussian distribution [33]. Pixels with too  
123 few observations were discarded, and a control quality check was applied to  
124 remove outliers [34].

125 Remotely sensed CHL may have missing data because of cloud coverage.  
126 The cloud variability in the Red Sea follows a seasonal cycle. Figure 1(c)  
127 shows that the cloud coverage is particularly pronounced during summers  
128 because of the monsoon and it is sparse during winters. The cloud coverage  
129 is, however, not homogenous over the Red Sea. It is much more pronounced  
130 in the south (figure 1(b)). In this region, almost no data are available during  
131 summers.

132 The dataset is separated into a training set, composed of the first seven  
133 years of data, and a validation set, containing the remaining three years.  
134 The computations for the data filling and for the covariance model selection  
135 and estimation are based only on the training dataset. The testing dataset

<sub>136</sub> is used to validate the model predictions outside the training period.

### <sub>137</sub> 3. Data filling and modeling

#### <sub>138</sub> 3.1. Data filling

<sub>139</sub> The DINEOF (Data Interpolating Empirical Orthogonal Function) is an  
<sub>140</sub> EOF-based, recursive method for the reconstruction of data matrices with  
<sub>141</sub> missing values [28, 29]. It estimates the values of the missing data by suc-  
<sub>142</sub> cessive singular values decompositions (SVD) of a given data matrix and  
<sub>143</sub> truncated reconstructions. The advantage of this method is that it does not  
<sub>144</sub> require any a priori information about the data. It has been successfully used  
<sub>145</sub> for reconstruction of incomplete chlorophyll datasets in different regions of  
<sub>146</sub> the ocean [35, 36, 37].

<sub>147</sub> Let  $\mathbf{X}$  be an  $m \times k$  centered data matrix with missing values initially filled  
<sub>148</sub> with 0s. Then, until the missing values have converged, the following steps  
<sub>149</sub> are repeated. An SVD is first applied to the data matrix:  $\mathbf{X} = \mathbf{U}\Sigma\mathbf{V}^T$ , with  
<sub>150</sub>  $\mathbf{U}$  an  $m \times m$  unitary matrix,  $\Sigma$  an  $m \times k$  diagonal matrix and  $\mathbf{V}$  a  $k \times k$  unitary  
<sub>151</sub> matrix. The missing values are then replaced by the truncated reconstruction  
<sub>152</sub> order  $n$  of the data matrix:  $\{\mathbf{X}\}_{i,j} = \{\mathbf{U}^{(n)}\Sigma^{(n)}(\mathbf{V}^{(n)})^T\}_{i,j}$ , for  $i, j$  indices  
<sub>153</sub> of the missing values, with  $\mathbf{U}^{(n)}$  the  $m \times n$  matrix composed of the  $n$  first  
<sub>154</sub> columns of  $\mathbf{U}$ ,  $\mathbf{V}^{(n)}$  the  $k \times n$  matrix composed of the  $n$  first columns of  
<sub>155</sub>  $\mathbf{V}$ , and  $\Sigma^{(n)}$  the  $n \times n$  diagonal matrix with the  $n$  largest eigenvalues on its  
<sub>156</sub> diagonal. It is assumed that the eigenvalues and eigenvector are sorted by  
<sub>157</sub> decreasing order of eigenvalues. In [29], the authors introduced the filtering

158 of the temporal covariance matrix as a way of reducing spurious oscillations  
159 that may appear when the data are sparsely sampled in time. This filtering  
160 is controlled by the parameter of the Laplacian filter and the number of times  
161 the filter is applied.

162 The values of the DINEOF parameters are determined following the  
163 method outlined in [29]. The smoothing parameter of the Laplacian filter  
164 is set to 0.005. The number of modes in the truncation and the number of  
165 times the filter is applied are chosen following a cross-validation technique. A  
166 random subset of observed values is taken from  $X$  and assumed to be missing  
167 before the DINEOF is applied. The algorithm is then run with different num-  
168 bers of iterations (1, 3, 10, 30, 100) and orders of truncation (from 2 to 50).  
169 The set of parameters minimizing the RMS error over the cross-validation  
170 data is chosen as the best number of iterations and order of truncation. The  
171 approach of [38] is followed to select a cross-validation dataset. Instead of  
172 selecting it by sampling the dataset point by point, contiguous regions are set  
173 aside. These regions correspond to regions of missing data from the original  
174 dataset and are selected randomly until 3% of the data have been extracted.

175 *3.2. Geostatiscal modeling*

The chlorophyll concentration data are modeled as a space-time random Gaussian process:

$$Z(\mathbf{s}; t), (\mathbf{s}; t) \in \mathbb{R}^2 \times \mathbb{R}. \quad (1)$$

<sub>176</sub> Such a process is entirely characterized by the mean function  $\mu(\mathbf{s}; t) =$   
<sub>177</sub>  $E(Z(\mathbf{s}; t))$ , and the covariance function  $K(\mathbf{s}, \mathbf{r}; t, q) = \text{cov}(Z(\mathbf{s}; t), Z(\mathbf{r}; q))$ .  
<sub>178</sub> The process is further assumed to be stationary, such that the covariance  
<sub>179</sub> function can be written as  $K(\mathbf{s}, \mathbf{r}; t, q) = \text{cov}(Z(\mathbf{s}; t), Z(\mathbf{r}; q))$ .

Covariance functions should be positive-definite, meaning that for any ensemble of space-time coordinates,  $\{(\mathbf{s}_i; t_i)\}_{i=1,\dots,k}$ , and real coefficients,  $d_1, \dots, d_k$ , the following condition should hold:

$$\sum_{i=1}^k \sum_{j=1}^k d_i d_j C(\mathbf{s}_i - \mathbf{s}_j; t_i - t_j) \geq 0. \quad (2)$$

<sub>180</sub> One way to enforce positive-definiteness in practice is to assume that the  
<sub>181</sub> covariance function belongs to a parametric family of covariance functions,  
<sub>182</sub> denoted by  $C(\mathbf{h}; u|\boldsymbol{\theta})$ , where  $\boldsymbol{\theta}$  is the vector of parameters to be estimated  
<sub>183</sub> [39].

### <sub>184</sub> 3.3. Mean function estimation and anomalies

<sub>185</sub> The data exhibit a distinct seasonal signal accounting for roughly 50%  
<sub>186</sub> of the total variability as can be seen in Figure 2. It is therefore reasonable  
<sub>187</sub> to model it in the mean function  $\mu(\mathbf{s}; t)$ . Once the data are filled using the  
<sub>188</sub> DINEOF algorithm, the seasonal signal is estimated for each week with the  
<sub>189</sub> weekly average computed from the training data.

<sub>190</sub> The anomalies are then computed by subtracting the weekly averages  
<sub>191</sub> from the data. The CHL anomalies of the first 8 weeks of data are displayed  
<sub>192</sub> in Figure 3. There are large regions of similar colors indicating spatial corre-

193 lations. Moreover, two maps adjacent in time display similar patterns, sug-  
194 gesting that the dataset is also correlated in time and justifying the use of  
195 space-time covariance functions to model CHL anomalies. From the anomaly  
196 time-series at three locations, plotted in Figure 4, one may conclude that the  
197 data have been successfully detrended. Since no specific pattern appears, the  
198 assumption of stationarity for the covariance model seems to be appropriate.

199 *3.4. Construction and fitting of the covariance model*

200 For space-time processes, families of covariance functions are typically  
201 built by combining known spatial and temporal covariance models. Sep-  
202 arable models are the simplest example of this approach, taking the form  
203  $C(\mathbf{h}; u|\boldsymbol{\theta}) = C^1(\mathbf{h}|\boldsymbol{\theta}_1)C^2(t|\boldsymbol{\theta}_2)$ , where  $C^1$  and  $C^2$  are respectively pure spa-  
204 tial and temporal covariance models parametrized by  $\boldsymbol{\theta}_1$  and  $\boldsymbol{\theta}_2$ . However,  
205 realizations from such families of covariance models were shown to exhibit  
206 non-physical behaviors [1, 27]. Research has been conducted to develop meth-  
207 ods for constructing non-separable covariance functions to overcome these  
208 limitations [27, 1].

In this work, families of space-time covariance functions are constructed following an approach proposed by Gneiting [1]. Given any completely monotonous function,  $\varphi(t), t \geq 0$ , and any positive function with a completely monotonous derivative,  $\psi(t), t \geq 0$ , the following space-time covari-

ance function is valid (i.e., positive-definite):

$$C(\mathbf{h}; u) = \frac{\sigma^2}{\psi(|u|^2)} \varphi \left( \frac{||\mathbf{h}||^2}{\psi(|u|^2)} \right), (\mathbf{h}, u) \in \mathbb{R}^2 \times \mathbb{R}. \quad (3)$$

- 209 In general,  $\varphi$  is defined as a temporal covariance function, and  $\psi$  as a spatial  
 210 covariance model, e.g., the exponential model or the Matérn model. Tables 1  
 211 and 2 present some examples of these functions that were introduced in [1].  
 212 To obtain more specialized families, a nugget term can be added.

The typical approach for estimating the parameters of the covariance function is by fitting it to the empirical space-time covariance matrix [39].

The empirical covariance function is first estimated using the formula:

$$2\hat{C}(\mathbf{h}(l); u) = \frac{1}{|N(\mathbf{h}(l); u)|} \sum_{(i,j,t,t') \in N(\mathbf{h}(l); u)} a(\mathbf{s}_i; t)a(\mathbf{s}_j; t'), \quad (4)$$

with  $N(\mathbf{h}(l); u) = \{(i, j, t, t') \text{ such that } : |t - t'| = u, \mathbf{s}_i - \mathbf{s}_j \in \text{Bin}(\mathbf{h}(l))\}$ , where  $\mathbb{R}^2$  has been divided into a finite number of bins,  $\text{Bin}(\mathbf{h}(l))$ , each of which has a representative lag,  $\mathbf{h}(l)$ . Here, bins are considered of the form  $\mathbf{k} \in \text{Bin}(\mathbf{h}(l))$  if  $(l - 1) * \delta h \leq ||\mathbf{k}|| \leq l\delta h, l = 0, \dots, L$ , with  $\delta$  a fixed interval size and  $L$  the number of retained intervals. The parameters are then estimated by fitting the covariance model using the Weighted Least Squares (WLS) approach. This method is popular due to its simplicity and proven efficiency [40]. In practice, the estimator of  $\boldsymbol{\theta}$  minimizes the weighted sum of the squares of

errors:

$$S(\boldsymbol{\theta}) = \sum_l \sum_{u=0}^3 w_{l,u} \left( \hat{C}(\mathbf{h}(l); u) - C(\mathbf{h}(l); u|\boldsymbol{\theta}) \right)^2, \quad (5)$$

213 with weights  $w_{l,u} = \frac{|N(\mathbf{h}(l); u)|}{(1-C(\mathbf{h}(l); u|\boldsymbol{\theta}))^2}$ .  $\phi$  and  $\psi$  are chosen from the functions  
214 presented in Tables 1 and 2, by separately fitting the covariance function in  
215 time and space for all possible choices of  $\phi$  and  $\psi$ . The choice of the functions  
216 is then made based on the quality of the fit.

217 **4. State-space formulation and filtering**

218 So far, the CHL data are modeled as  $Z(\mathbf{s}; t) = \mu(\mathbf{s}; t) + a(\mathbf{s}; t)$ , where  $\mu$   
219 is a deterministic function representing the seasonality and  $a$  is a zero-mean  
220 space-time Gaussian process representing the weekly anomalies, with the  
221 covariance function estimated as described in the previous section. However,  
222 a state-space model needs to be formulated to use the Kalman filter. This  
223 is done by identifying a CHL underlying process,  $Y(\mathbf{s}; t)$ , distinct from the  
224 data process  $Z(\mathbf{s}; t)$ . These processes are then discretized in space and time  
225 to obtain the state-space formulation, on which a reduced-order variant of  
226 the Kalman filter is applied.

227 4.1. State-space modeling

To derive the state-space model, we start from the following model:

$$Y(\mathbf{s}; t) = \mu(\mathbf{s}; t) + a(\mathbf{s}; t), \text{ where } \mathbf{s}, t \text{ are observable,} \quad (6)$$

$$Z(\mathbf{s}; t) = Y(\mathbf{s}; t) + \varepsilon(\mathbf{s}; t), \text{ where } \mathbf{s}, t \text{ are observed.} \quad (7)$$

228 In this formulation,  $\varepsilon(\mathbf{s}; t)$  is a white noise process representing the mea-  
229 surement errors,  $a(\mathbf{s}; t)$  is the anomaly process with the covariance matrix  
230 estimated as described in the previous section,  $\mu(\mathbf{s}; t)$  is the seasonal mean  
231 function,  $Y(\mathbf{s}; t)$  is the underlying CHL log-concentration process, and  $Z(\mathbf{s}; t)$   
232 the data process.

To build a state-space model for filtering and forecasting the spatial variability of CHL in time, once the covariance function is estimated, we resort to Kriging. In classical geostatistics, Kriging interpolates available observations to provide the best linear unbiased estimate of a spatial process at unobserved locations [39]. This technique easily generalizes to space-time Gaussian processes, particularly for forecasting. The equations are obtained by conditioning the anomalies at time  $t$  by the anomalies at time  $t - 1$ . Using the vectorial notations  $\mathbf{a}_t = \{a(\mathbf{s}_{i,t})\}_{i=1,\dots,k}$ , where  $\mathbf{s}_{i,t}$  is the spatial location of the  $i$ -th observation at time  $t$ ,  $[\mathbf{a}_t^T, \mathbf{a}_{t-1}^T]^T$  is a Gaussian vector. Therefore,

by conditioning

$$\mathbf{a}_t | \mathbf{a}_{t-1} \sim N(\mathbf{M}\mathbf{a}_{t-1}, \mathbf{Q}), \quad (8)$$

$$(9)$$

where

$$\mathbf{M} = \mathbf{C}_1 \mathbf{C}_0^{-1}, \quad (10)$$

$$\mathbf{Q} = \mathbf{C}_0 - \mathbf{C}_1 \mathbf{C}_0^{-1} \mathbf{C}_1, \quad (11)$$

with  $\mathbf{C}_0 = E[\mathbf{a}_t \mathbf{a}_t^T] = \{C(\mathbf{s}_{i,t} - \mathbf{s}_{j,t}), 0\}_{i,j=1,\dots,k}$  and  $\mathbf{C}_1 = E[\mathbf{a}_t \mathbf{a}_{t-1}^T] = \{C(\mathbf{s}_{i,t} - \mathbf{s}_{j,t-1}), 1\}_{i,j=1,\dots,k}$ , the following recursive state-space model can be derived:

$$\mathbf{y}_t = \boldsymbol{\mu}_t + \mathbf{M}(\mathbf{y}_{t-1} - \boldsymbol{\mu}_{t-1}) + \boldsymbol{\eta}_t, \quad (12)$$

$$\mathbf{z}_t = \mathbf{H}_t \mathbf{y}_t + \boldsymbol{\varepsilon}_t, \quad (13)$$

and  $\mathbf{y}_t = \mathbf{a}_t + \boldsymbol{\mu}_t$  is the CHL process discretized in space,  $\boldsymbol{\mu}_t$  is the vector seasonal component,  $H_t$  is the observation operator that returns the CHL concentration at the observed location,  $\boldsymbol{\eta}_t$  represents the model error, and  $\boldsymbol{\varepsilon}_t$  represents the measurement error.  $\mathbf{y}_t$  is a vector of fixed size that represents the whole model domain (Red Sea), whereas  $\mathbf{z}_t$  has a variable size equal to the number of available observations at time  $t$ .

The model can be reformulated such that only the anomalies are filtered

and  $a_t$  is the state vector:

$$\mathbf{a}_t = \mathbf{M}\mathbf{a}_{t-1} + \boldsymbol{\eta}_t, \quad (14)$$

$$\mathbf{z}_t = \mathbf{y}_t - \mathbf{H}_t \mathbf{s}_t = \mathbf{H}_t \mathbf{a}_t + \boldsymbol{\varepsilon}_t. \quad (15)$$

This system is equivalent to the preceding one, but is more practical to implement. Given  $\boldsymbol{\eta}_t \sim N(0, \mathbf{Q})$  and assuming  $\boldsymbol{\varepsilon}_t \sim N(0, \sigma_{\text{obs}}^2 \mathbf{I})$ , independent and identically distributed (i.i.d.), with  $\sigma_{\text{obs}}^2$  a fixed constant that is tuned by minimizing the RMS prediction error.

#### 4.2. Low-rank Kalman filter

Two issues need to be tackled before using the time evolution model for predicting the anomalies. First, the consequent amount of missing data in the CHL satellite observations makes it difficult to obtain frequent initial CHL concentrations to integrate the model forward in time for forecasting, and second, the model is linear with eigenvalues smaller than one<sup>1</sup> in absolute values, making it inappropriate for long-term predictions. The latter means that CHL model forecasts would decrease exponentially in time, quickly becoming close to 0. The Kalman filter solves both problems by recursively assimilating the observations and providing an optimal estimate, in the mean-square sense, of the anomalies to start a new forecast cycle. The filter operates in two steps to compute the best linear estimate of the state

---

<sup>1</sup>This was verified numerically in the present work. We are not aware of a general result.

255 of a linear dynamical model given past observations [41].

- *Forecast step:* Starting from the best available estimate of the state,  $a_{t-1}^a$ , and the associated error covariance matrix,  $P_{t-1}^a$ , at a given time,  $t - 1$ , the forecast state,  $a_t^f$ , and its error covariance matrix,  $P_t^f$ , are obtained by integrating the model forward to the time of the next available observation:

$$\mathbf{a}_t^f = \mathbf{M}\mathbf{a}_{t-1}^a, \quad (16)$$

$$\mathbf{P}_t^f = \mathbf{M}\mathbf{P}_{t-1}^a\mathbf{M}^T + \mathbf{Q}. \quad (17)$$

- *Update step:* Once a new observation,  $\mathbf{z}_t$ , is available, the forecast state,  $\mathbf{a}_t^f$ , and its error covariance,  $\mathbf{P}_t^f$ , are updated to their analysis counterparts,  $\mathbf{a}_t^a$ , and,  $\mathbf{P}_t^a$ , as:

$$\mathbf{a}_t^a = \mathbf{a}_t^f + \mathbf{K}_t(\mathbf{z}_t - \mathbf{H}_t\mathbf{a}_t^f), \quad (18)$$

$$\mathbf{P}_t^a = (\mathbf{I} - \mathbf{K}_t\mathbf{H}_t)\mathbf{P}_t^f, \quad (19)$$

$$\mathbf{K}_t = \mathbf{P}_t^f \mathbf{H}_t^T (\mathbf{H}_t \mathbf{P}_t^f \mathbf{H}_t^T + \sigma_{\text{obs}}^2 \mathbf{I})^{-1} \quad (20)$$

256 where  $K_t$  is known as the Kalman gain.

The inversion in the computation of the Kalman gain is computationally quite demanding because the number of observations is as large as the size of the state. To speed up this computation, the reduced-order Kalman filter [31]

is used. This filter approximates the covariance matrices as  $\mathbf{P}_t^f = \mathbf{L}\mathbf{U}_t^f\mathbf{L}^T$  and  $\mathbf{P}_t^a = \mathbf{L}\mathbf{U}_t^a\mathbf{L}^T$ , where  $\mathbf{L}$  is a  $n \times r$  matrix whose columns are the  $r$  leading EOFs (computed here with the DINEOF algorithm), and  $\mathbf{U}_t^f$  and  $\mathbf{U}_t^a$  are  $r \times r$  matrices. The inversion is then applied on the  $r \times r$  matrices  $\mathbf{U}_t^f$  and  $\mathbf{U}_t^a$ , which considerably reduces the computational burden since  $r \ll n$  in practice [31]. Defining

$$\mathbf{V} = \mathbf{L}^T \mathbf{Q} \mathbf{L}, \quad (21)$$

$$\mathbf{W} = \mathbf{L}^T \mathbf{M} \mathbf{L}, \quad (22)$$

which respectively represent the projection of the model dynamics and the model error on the leading EOFs, the reduced-order Kalman filter equations for the forecast and update of the covariance matrix can be simplified by only updating  $\mathbf{U}_t^f$  and  $\mathbf{U}_t^a$  as follows:

$$\mathbf{U}_t^f = \mathbf{W} \mathbf{U}_{t-1}^a \mathbf{W}^T + \mathbf{V}, \quad (23)$$

$$(\mathbf{U}_t^a)^{-1} = (\mathbf{U}_t^f)^{-1} + \mathbf{L}^T \mathbf{H}^T \mathbf{R}^{-1} \mathbf{H} \mathbf{L}. \quad (24)$$

The Kalman gain can be then computed as

$$\mathbf{K}_t = \sigma_{\text{obs}}^{-2} \mathbf{L} \mathbf{U}_t^a \mathbf{L}^T \mathbf{H}^T. \quad (25)$$

<sup>257</sup> Otherwise, the forecast and update steps are identical to those of the Kalman  
<sup>258</sup> filter.

259    4.3. Low-rank Kalman filter with colored noise

260    One issue is that the model noise is correlated in time, whereas the  
 261    Kalman filter assumes it to be white. Indeed,  $\boldsymbol{\eta}_t = \mathbf{a}_t - \mathbf{M}\mathbf{a}_{t-1}$  is a space-  
 262    time Gaussian process. One can show that  $[\boldsymbol{\eta}_t^T, \boldsymbol{\eta}_{t-1}^T]^T$  is a Gaussian vector  
 263    of mean 0 and covariance matrix  $\begin{bmatrix} \mathbf{A}_0 & \mathbf{A}_1 \\ \mathbf{A}_1^T & \mathbf{A}_0 \end{bmatrix}$ , with  $\mathbf{A}_0 = \mathbf{C}_0 - \mathbf{C}_1\mathbf{C}_0^{-1}\mathbf{C}_1^T$ ,  
 264     $\mathbf{A}_1 = \mathbf{C}_1\mathbf{C}_0^{-1}\mathbf{C}_1\mathbf{C}_0^{-1}\mathbf{C}_1 - \mathbf{C}_2\mathbf{C}_0^{-1}\mathbf{C}_1$ . We can then predict  $\boldsymbol{\eta}_t$  by condition-  
 265    ing on the previous equation and derive the model  $\boldsymbol{\eta}_t = \mathbf{N}\boldsymbol{\eta}_{t-1} + \boldsymbol{\xi}_t$ , with  
 266     $\mathbf{N} = \mathbf{A}_1\mathbf{A}_0^{-1}$  and  $\boldsymbol{\xi}_t \sim \mathcal{N}(0, \mathbf{A}_0 - \mathbf{A}_1\mathbf{A}_0^{-1}\mathbf{A}_1)$ .

267    To take into account this correlation in the Kalman filter, the state-space  
 268    system is enlarged by including  $\boldsymbol{\eta}_t$  in the state vector [32], so the new state-  
 269    space system becomes:

$$\hat{\mathbf{a}}_t = \boldsymbol{\Gamma}\hat{\mathbf{a}}_{t-1} + \hat{\boldsymbol{\eta}}_t, \quad (26)$$

$$\mathbf{y}_t = \hat{\mathbf{H}}\hat{\mathbf{a}}_{t-1} + \boldsymbol{\varepsilon}_t, \quad (27)$$

270    with  $\hat{\mathbf{a}}_t = [\mathbf{a}_t^T, \boldsymbol{\eta}_t^T]^T$ ,  $\boldsymbol{\Gamma} = \begin{bmatrix} \mathbf{M} & \mathbf{I} \\ \mathbf{0} & \mathbf{N} \end{bmatrix}$ ,  $\hat{\boldsymbol{\eta}}_t = [\mathbf{0}, \boldsymbol{\xi}_t^T]^T$  and  $\hat{\mathbf{H}}_t = [\mathbf{H}_t, \mathbf{0}]$ , on  
 271    which the Kalman filter is applied. Of course,  $\boldsymbol{\xi}_t$  is again colored, but the  
 272    procedure can be iterated until no significant correlation is left to exploit in  
 273    the noise.

274 **5. Results**

275 Figure 5 summarizes the workflow of the experimental setup. The chloro-  
276 phyll concentration is the sum of a seasonal component and an anomaly. The  
277 former is estimated using data filled by the DINEOF algorithm over a learn-  
278 ing period (first seven years of CHL data). The anomalies are assumed to be  
279 a stationary Gaussian process whose space-time covariance matrix has been  
280 estimated. A reduced-order Kalman filter is then applied to perform one-  
281 step ahead predictions. These are compared with the observations over the  
282 remaining three years of validation data to validate the system’s performance  
283 and results. The results of the DINEOF data filling algorithm are presented  
284 before a covariance model for the data is chosen and fitted. The results of  
285 the filtering are finally examined and analyzed.

286 *5.1. DINEOF Analysis*

287 Figure 6 plots the (RMS) error for the cross-validation dataset and for  
288 different values of the number of smoothing iterations and EOFs. A lower  
289 RMS error is achieved with 24 EOFs and 30 iterations.

290 Figures 7(a) and 7(c) show the first spatial and associated temporal modes  
291 resulting from the DINEOF analysis. Figure 7(c) exhibits a regular peak dur-  
292 ing winters and a secondary peak of varying size during summers. These can  
293 be associated with the winter bloom and the secondary summer bloom, both  
294 described in [9]. Figure 7(a) shows that the bloom is relatively homoge-  
295 neous over the Red Sea, except in the southwest corner where the variation

296 is strongly pronounced.

297 Figures 7(b) and 7(d) plot the second spatial and temporal modes re-  
298 spectively. The spatial mode shows a north-south contrast. The time series  
299 displays a large peak around summer 2000, which corresponds to a large  
300 positive anomaly taking place during this period (Figure 7(e)). Except for  
301 the first two modes, all other modes tend to explain a local feature of the  
302 data. This makes the interpretation of the EOF analysis difficult and the  
303 convergence of the spectrum very slow as shown in Figure 7(f).

304 *5.2. Covariance model estimation*

305 Figure 8a plots the empirical space covariance function  $\hat{C}(\|\mathbf{h}\|, 0)$ . Among  
306 the functions in Table 1, the ones that best fit this curve are selected. A WLS  
307 minimization is applied to fit the parameters  $c$ ,  $\gamma$  and  $\nu$ , as well as a the noise  
308 variance  $\sigma_{\text{spatial}}^2$  while ignoring the observations for  $\|\mathbf{h}\| = 0$ . This allows a  
309 nugget effect to be taken account. The results are superimposed on Figure  
310 8a. The Matérn model clearly fits poorly. By examining the residuals, we  
311 can obviously eliminate  $\varphi_4$ . Finally, between the two remaining candidates,  
312 both of which appear to be plausible,  $\varphi_1$  is chosen since it involves fewer  
313 degrees of freedom.

314 Figure 8b plots the empirical time covariance function  $\hat{C}(0; u)$ .  $\psi_2$  defines  
315 a time covariance that decreases very slowly. The covariance with  $\psi_3$  con-  
316 verges to a constant as the time-lag goes to infinity, which is not realistic in  
317 the case of anomalies. The first-time covariance model is filled with a nugget

318 effect for  $u = 0$ , using WLS, and the result is displayed in Figure 8b, where  
 319 a nugget effect can be clearly identified.

To build the space-time covariance function, we use the reparametrization in [1] (see example 1), along with the building blocks in time and space that have been derived in Section 3. A time-only nugget effect is also added, leading to the following covariance model:

$$C(\mathbf{h}; u) = \begin{cases} \sigma^2 \exp(-c||\mathbf{h}||^{2\gamma}) & , \text{ if } u = 0, \\ \frac{\sigma^2 \tau}{a|u|^{2\alpha} + 1} \exp\left(-\frac{c||\mathbf{h}||^{2\gamma}}{(a|u|^{2\alpha} + 1)^{\beta\gamma}}\right) & , \text{ otherwise,} \end{cases} \quad (28)$$

320 with  $0 < \tau \leq 1$ .

321 The initial values for the WLS fitting are given by the results of the  
 322 preceding purely spatial and purely temporal regressions. The results are  
 323 shown in Table 3. A contour plot of the resulting covariance function is  
 324 shown in Figure 9. Our fitted correlation model has level curves very close  
 325 to those of the empirical covariance model.

### 326 5.3. Filtering

327  $\sigma_{\text{obs}}^2$  is first determined empirically by trial and error, chosen as the value  
 328 that leads to the minimal averaged RMS error (RMSE).  $\sigma_{\text{obs}}^2 = 1$  is found to  
 329 be a reasonable choice.

330 The anomaly model with a space-time covariance function helps to de-  
 331 crease the RMS error by nearly 11% over the test period. As Figure 10a indi-  
 332 cates, the improvement is most noticeable during periods with large anomalies.

333 The boxplot in Figure 10b shows that the variability of errors is reduced with  
334 the proposed filtering approach. Since the modeling is purely statistical and  
335 is not based on physical quantities, it was not able to anticipate the start  
336 of a bloom. However, it successfully extrapolates the current estimate of  
337 the anomaly in space and time and improves the prediction compared to the  
338 seasonal component.

339 In Figures 11 (a to c), the predictions of the model are compared with  
340 the observations and the seasonal signal for a fall week in 2006. The model  
341 captures some differences with the seasonal regime, such as a larger northern  
342 region with a low CHL concentration extending south below 22°N, and a  
343 more intense bloom in the south. The usual seasonal dynamic is altered with  
344 an extension of the stratified regime in the north and a larger intrusion of  
345 the nutrient rich waters of the Gulf of Aden [9].

346 Figure 11 (d to f) shows a similar comparison for a winter week in 2006.  
347 The model captures a weak winter bloom in the northern half of the Red  
348 Sea. A similar pattern has been described in [7] for winter 1999, and it  
349 seems to be a common feature of El Niña years. The model also successfully  
350 captures a high CHL concentration in the south, and a lower than usual CHL  
351 concentration in the center [9].

352 Figure 12 plots the error variance as predicted by the filter and the actual  
353 prediction RMSE, computed as the difference between the model forecast  
354 and data, averaged over the three year training period at every point of the  
355 grid. The predicted RMSE corresponds to the diagonals of the prediction

356 covariance matrix as estimated by the reduced-order Kalman filter. The  
357 prediction RMSE is computed from the error between the model prediction  
358 and the observation. We can see that both values are close in the northern  
359 half of the Red Sea, but the RMSE is much larger in the south. This is caused  
360 by the lack of data and the fact that the dynamics in the South is different  
361 from that in the north, making the process nonstationary. Indeed, as shown  
362 in Figure 3, the anomalies in the south clearly exhibit smaller spatial and  
363 temporal correlation length scales compared with those in the north.

364 One way to evaluate the filter’s behavior is to examine the distribution of  
365 the innovations and the increments [42]. The innovation corresponds to the  
366 difference:  $\mathbf{y}_t - \mathbf{H}_t \mathbf{a}_t^f$ , whereas the increment corresponds to the difference:  
367  $\mathbf{H}_t (\mathbf{a}_t^a - \mathbf{a}_t^f)$ . Figure 13(c) shows that the innovation seems to be approxi-  
368 mately normally distributed, as expected for a properly tuned Kalman filter.  
369 Figure 13(b) shows that the averaged innovation size decreases to a value  
370 close to zero as the filter assimilates the data over time, which also suggests  
371 that the filter is properly working. However, Figure 13(a) indicates that the  
372 increments of the filter tend to be positively biased in some regions. This indi-  
373 cates that the model tends to underestimate the amount of CHL, which may  
374 be associated with the statistical model’s inability to forecast CHL blooms.

375 Figures 14 (a), (b) and (c) plot the spatially averaged correlations between  
376 the observations and the seasonal predictions, the Kalman filter forecasts and  
377 the analyses, respectively. Compared with the seasonal correlations, we can  
378 see that the model improves the prediction skill over the entire Red Sea,

379 particularly in its central part, with correlations ranging between 0.28 and  
380 0.92. The filter further improves the prediction-data correlation, over the  
381 entire domain with correlations up to 0.96.

382 **6. Discussion**

383 Here, we considered the filtering problem of satellite-derived chlorophyll  
384 (CHL) concentration in the Red Sea using a data-driven approach in which  
385 the CHL spatio-temporal evolution is modeled as a space-time Gaussian pro-  
386 cess. The DINEOF data-filling algorithm [38] was applied to compute an  
387 estimate of the seasonal signal in the data, which was used as the mean  
388 function of the process. To model the residual anomalies, the method pro-  
389 posed by Gneiting in [1] was applied to construct an appropriate family of  
390 covariance functions.

391 From the anomalies, a family of covariance functions is constructed and  
392 then fitted to the data. Based on a space-time Kriging formulation, a linear  
393 model was derived to capture changes in the chlorophyll concentration. A  
394 reduced-order variant of the Kalman filter was then applied to forecast and  
395 filter the CHL concentration. The results of our experiments suggest that  
396 the proposed system works reasonably well, reducing the RMS error in CHL  
397 concentration prediction by about 11% as compared with the seasonal mean.

398 The proposed approach is not difficult to apply, but requires some coding

efforts. For the DINEOF, a standalone package is freely available online<sup>2</sup>. To the authors' knowledge, there is currently no R package or library for fitting custom space-time covariance models. The reduced-order Kalman filter is straightforward to implement.

The proposed method requires the estimation of a very few parameters, which may prevent overfitting. Another advantage is its stability, as the anomaly prediction cannot grow over time. One problem with Kriging is that in the case of missing data over large areas, its prediction will be the mean function [30]. In this study, this problem is alleviated using the Kalman filter which always provides a prediction over the whole domain.

Forecasting and filtering CHL concentration in the Red Sea using remotely sensed data is challenging. Because of cloud coverage over the southern Red Sea during the summer, large areas remain unobserved. Moreover, the Red Sea is a heterogenous environment with different ecological and physical dynamics from north to south. We can therefore expect the anomalies to be non-stationary. Fortunately, Kriging is robust to non-stationarity [30]. A problem with this correlation-based modeling is its inability to forecast CHL blooms. Our linear model tends to underestimate the amount of CHL concentration and to miss blooms in some regions, before the Kalman filter corrects this error when new observations are available.

The proposed method can be further generalized by considering more

---

<sup>2</sup><http://modb.oce.ulg.ac.be/mediawiki/index.php/DINEOF>

420 sophisticated mean functions. For example, one can use a linear model with  
421 additional covariates, such as sea surface temperature, sea surface height,  
422 thermocline depth or wind speed. Another way to improve the model is to  
423 use non-stationary covariance functions, or to use non-Gaussian models to  
424 predict blooms. These tasks will be considered in future studies.

425 **Acknowledgment**

426 The research reported in this publication was supported by King Abdullah  
427 University of Science and Technology (KAUST).

428 **7. Bibliography**

429 [1] T. Gneiting, Nonseparable, Stationary Covariance Functions for Space-  
430 Time Data, *Journal of the American Statistical Association* 97 (458)  
431 (2002) 590–600. doi:10.1198/016214502760047113.

432 [2] F. Yao, I. Hoteit, L. J. Pratt, A. S. Bower, P. Zhai, A. Köhl, G. Gopalakrishnan,  
433 Seasonal overturning circulation in the Red Sea: part 1. Model validation and summer circulation,  
434 *Journal of Geophysical Research: Oceans* (2014) 2238–2262doi:10.1002/2013JC009331.Key.

436 [3] S. Heileman, N. Mistafa, III-6 Red Sea: LME# 33, lme.noaa.gov.  
437 URL [http://www.lme.noaa.gov/lmeweb/LME\\_Report/lme\\_33.pdf](http://www.lme.noaa.gov/lmeweb/LME_Report/lme_33.pdf)

438 [4] D. E. Raitsos, I. Hoteit, P. K. Prihartato, T. Chronis, G. Triantafyllou,

- 439 Y. Abualnaja, Abrupt warming of the Red Sea, Geophysical Research  
440 Letters 38. doi:10.1029/2011GL047984.
- 441 [5] A. Lo-Yat, S. D. Simpson, M. Meekan, D. Lecchini, E. Martinez,  
442 R. Galzin, Extreme climatic events reduce ocean productivity and  
443 larval supply in a tropical reef ecosystem, Global Change Biology 17  
444 (2011) 1695–1702. doi:10.1111/j.1365-2486.2010.02355.x.  
445 URL <http://onlinelibrary.wiley.com/doi/10.1111/j.1365-2486.2010.02355.x/full>
- 446 [6] K. H. Mann, J. R. N. Lazier, Dynamics of Marine Ecosystems, 2006.
- 447 [7] J. Acker, G. Leptoukh, S. Shen, T. Zhu, S. Kempler, Remotely-sensed  
448 chlorophyll a observations of the northern Red Sea indicate seasonal  
449 variability and influence of coastal reefs, Journal of Marine Systems 69  
450 (2008) 191–204. doi:10.1016/j.jmarsys.2005.12.006.  
451 URL <http://linkinghub.elsevier.com/retrieve/pii/S0924796307000589>
- 452 [8] G. Triantafyllou, F. Yao, G. Petihakis, K. Tsiaras, D. E. Raitsos,  
453 I. Hoteit, Exploring the Red Sea seasonal ecosystem functioning using a  
454 three-dimensional biophysical model, Journal of Geophysical Research:  
455 Oceans (2014) 1791–1811doi:10.1002/2013JC009641.Received.
- 456 [9] D. E. Raitsos, Y. Pradhan, R. J. W. Brewin, G. Stenchikov, I. Hoteit,  
457 Remote Sensing the Phytoplankton Seasonal Succession of the Red Sea,  
458 PLoS ONE 8 (6). doi:10.1371/journal.pone.0064909.

- 459 [10] R. J. W. Brewin, D. E. Raitsos, Y. Pradhan, I. Hoteit, Comparison  
460 of chlorophyll in the Red Sea derived from MODIS-Aqua and in  
461 vivo fluorescence, *Remote Sensing of Environment* 136 (2013) 218–224.  
462 doi:10.1016/j.rse.2013.04.018.
- 463 URL <http://dx.doi.org/10.1016/j.rse.2013.04.018>
- 464 [11] C. R. McClain, A Decade of Satellite Ocean Color Ob-  
465 servations, *Annual Review of Marine Science* (2009) 19–  
466 42doi:10.1146/annurev.marine.010908.163650.
- 467 [12] T. R. Anderson, Plankton functional type modelling: running  
468 before we can walk?, *Journal of Plankton Research* 27 (11).  
469 doi:10.1093/plankt/fbi076.
- 470 URL <http://www.plankt.oxfordjournals.org/cgi/doi/10.1093/plankt/fbi076>
- 471 [13] J. Baretta, W. Ebenhöh, P. Ruardij, The European regional seas  
472 ecosystem model, a complex marine ecosystem model, *Netherlands*  
473 *Journal of Sea Research* 33 (3/4) (1995) 233–246. doi:10.1016/0077-  
474 7579(95)90047-0.
- 475 URL <http://linkinghub.elsevier.com/retrieve/pii/0077757995900470>
- 476 [14] G. Petihakis, G. Triantafyllou, I. J. Allen, I. Hoteit, C. Dounas, Mod-  
477 elling the spatial and temporal variability of the Cretan Sea ecosystem,  
478 *Journal of Marine Systems* 36 (2002) 173–196.
- 479 [15] J. C. Blackford, P. J. Radford, A Structure and methodology for ma-

- 480 marine ecosystem modelling, Netherlands Journal of Sea Research 33 (3/4)  
481 (1995) 247–260.

482 [16] J. H. W. Lee, Y. Huang, M. Dickman, A. W. Jayawardena, Neural  
483 network modelling of coastal algal blooms 159 (2003) 179–201.

484 [17] F. Recknagel, Artifical neural network approach for modelling and pre-  
485 diction of algal blooms, Ecological Modelling 96 (1997) 11–28.

486 [18] D. Raitsos, G. Korres, G. Triantafyllou, G. Petihakis, M. Pantazi,  
487 K. Tsiaras, A. Pollani, Assessing chlorophyll variability in relation  
488 to the environmental regime in Pagasitikos Gulf, Greece, Journal of  
489 Marine Systems 94 (2012) S16–S22. doi:10.1016/j.jmarsys.2011.11.003.  
490 URL <http://linkinghub.elsevier.com/retrieve/pii/S0924796311002703>

491 [19] D. E. Raitsos, S. J. Lavender, Y. Pradhan, T. Tyrrell, P. C. Reid, M. Ed-  
492 wards, Coccolithophore bloom size variation in response to the regional  
493 environment of the subarctic North Atlantic, Limnology and Oceanog-  
494 raphy 51 (5) (2006) 2122–2130. doi:10.4319/lo.2006.51.5.2122.  
495 URL [http://www.aslo.org/lo/toc/vol\\_51/issue\\_5/2122.html](http://www.aslo.org/lo/toc/vol_51/issue_5/2122.html)

496 [20] P. C. Kyriakidis, A. G. Journel, Geostatistical Space-Time Models: A  
497 Review, Mathematical Geology 31 (6) (1999) 651–684.

498 [21] N. Cressie, C. K. Wikle, Statistics for Spatio-Temporal Data, 2011.

499 [22] M. S. Handcock, J. R. Wallis, An Approach to Statistical Spatial-

- 500        Temporal Modeling of Meteorological Fields, *Journal of the American*  
501        Statistical Association
- 502 [23] G. R. North, J. Wang, M. G. Genton, Correlation Models for  
503        Temperature Fields, *Journal of Climate* 24 (2011) 5850–5862.  
504        doi:10.1175/2011JCLI4199.1.
- 505 [24] M. E. Hohn, A. M. Liebhold, L. S. Gribko, Geostatistical Model for  
506        Forecasting Spatial Dynamics of Defoliation Caused by the Gypsy Moth  
507        (Lepidoptera: Lymantriidae), *Environmental Entomology* 22 (5) (1993)  
508        1066–1075.
- 509 [25] R. Lark, Towards soil geostatistics, *Spatial Statistics* 1 (2012) 92–99.
- 510 [26] A. Keaney, J. McKinley, C. Graham, M. Robinson, A. Ruffell, Spatial  
511        statistics to estimate peat thickness using airborne radiometric data,  
512        *Spatial Statistics* 5 (2013) 3–24.
- 513 [27] N. Cressie, H.-C. Huang, Classes of Nonseparable, Spatio-Temporal Sta-  
514        tionary Covariance Functions, *Journal of the American Statistical As-  
515        sociation* 94 (448) (1999) 1330–1340.
- 516 [28] J.-M. Beckers, M. Rixen, EOF Calculations and Data Filling from In-  
517        complete Oceanographic Datasets, *Journal of Atmospheric and Oceanic*  
518        *Technology* 20 (2003) 1839–1856.
- 519 [29] A. Alvera-Azcárate, A. Barth, D. Sirjacobs, J.-M. Beckers, Enhancing

- 520 temporal correlations in EOF expansions for the reconstruction of miss-  
521 ing data using DINEOF, Ocean Science 5 (2009) 475–485.
- 522 [30] P. Monestiez, L. Dubroca, E. Bonnin, J.-P. Durbec, C. Guinet, Geo-  
523 statistical modelling of spatial distribution of Balaenoptera physalus  
524 in the Northwestern Mediterranean Sea from sparse count data and  
525 heterogeneous observation efforts, Ecological Modelling 193 (3-4) (2006)  
526 615–628. doi:10.1016/j.ecolmodel.2005.08.042.  
527 URL <http://linkinghub.elsevier.com/retrieve/pii/S0304380005004436>
- 528 [31] I. Hoteit, D.-T. Pham, J. Blum, A simplified reduced order Kalman fil-  
529 tering and application to altimetric data assimilation in Tropical Pacific,  
530 Journal of Marine Systems 36 (2002) 101–127.
- 531 [32] C. K. Chui, G. Chen, Kalman Filtering, with Real-Time Applications,  
532 Fourth Edition, 2009.
- 533 [33] J. W. Campbell, The lognormal distribution as a model for bio-optical  
534 variability in the sea, Journal of Geophysical Research 100 (13) (1995)  
535 237–254.
- 536 [34] J. K. Willis, Interannual variability in upper ocean heat content, tem-  
537 perature, and thermosteric expansion on global scales, Journal of Geo-  
538 physical Research 109. doi:10.1029/2003JC002260.  
539 URL <http://doi.wiley.com/10.1029/2003JC002260>

- 540 [35] T. N. Miles, R. He, Temporal and spatial variability of Chl-a and SST  
541 on the South Atlantic Bight: Revisiting with cloud-free reconstructions  
542 of MODIS satellite imagery, *Continental Shelf Research* 30 (2010) 1951–  
543 1962. doi:10.1016/j.csr.2010.08.016.  
544 URL <http://dx.doi.org/10.1016/j.csr.2010.08.016>
- 545 [36] D. Sirjacobs, A. Alvera-Azcárate, A. Barth, G. Lacroix, Y. Park,  
546 B. Nechad, K. Ruddick, J.-M. Beckers, Cloud filling of ocean colour  
547 and sea surface temperature remote sensing products over the South-  
548 ern North Sea by the Data Interpolating Empirical Orthogonal  
549 Functions methodology, *Journal of Sea Research* 65 (2011) 114–130.  
550 doi:10.1016/j.seares.2010.08.002.  
551 URL <http://linkinghub.elsevier.com/retrieve/pii/S1385110110001036>
- 552 [37] J. N. Waite, F. J. Mueter, Spatial and temporal variability of  
553 chlorophyll-a concentrations in the coastal Gulf of Alaska, 1998–  
554 2011, using cloud-free reconstructions of SeaWiFS and MODIS-  
555 Aqua data, *Progress in Oceanography* 116 (2013) 179–192.  
556 doi:10.1016/j.pocean.2013.07.006.  
557 URL <http://dx.doi.org/10.1016/j.pocean.2013.07.006>
- 558 [38] J.-M. Beckers, A. Barth, A. Alvera-Azcárate, DINEOF reconstruction  
559 of clouded images including error maps-application to the Sea-Surface  
560 Temperature around Corsican Island, *Ocean Science* 2 (2006) 183–199.

- 561 doi:10.5194/osd-3-735-2006.
- 562 URL <http://www.ocean-sci-discuss.net/3/735/2006/>
- 563 [39] D. L. Zimmerman, M. Stein, Classical Geostatistical Methods,  
564 in: A. Gelfand, P. Diggle, M. Fuentes, P. Guttorp (Eds.),  
565 Handbook of Spatial Statistics, Vol. 20103158 of Chapman &  
566 Hall/CRC Handbooks of Modern Statistical Methods, 2010, pp. 29–44.  
567 doi:10.1201/9781420072884.
- 568 [40] D. L. Zimmerman, M. B. Zimmerman, A Comparison of Spatial Semi-  
569 variogram Estimators and Corresponding Ordinary Kriging Predictors,  
570 Technometrics 33 (1) (1991) 77–91.
- 571 [41] R. H. Shumway, S. S. David, Time Series Analysis and Its Applications,  
572 2011.
- 573 [42] P. Brasseur, Ocean data assimilation using sequential methods based  
574 on the Kalman filter, in: Ocean Weather Forecasting, Springer Edition,  
575 2006, pp. 271–316.

<sup>576</sup> **List of Tables**

<sup>577</sup> 1	Spatial building blocks proposed by [1] . . . . .	36
<sup>578</sup> 2	Temporal building blocks proposed by [1] . . . . .	37
<sup>579</sup> 3	Parameters of the covariance model given in equation (28) estimated by WLS. . . . .	38
<sup>580</sup>		

Table 1: Spatial building blocks proposed by [1]

Function	Parameters
$\varphi_1(t) = \exp(-ct^\gamma)$	$c > 0, 0 < \gamma \leq 1$
$\varphi_2(t) = (2^{\nu-1}\Gamma(\nu))^{-1}(ct^{1/2})^\nu K_\nu(ct^{1/2})$	$c > 0, \nu > 0$
$\varphi_3(t) = (1 + ct^\gamma)^{-\nu}$	$c > 0, 0 < \gamma \leq 1, \nu > 0$
$\varphi_4(t) = 2^\nu(\exp(ct^{1/2}) + \exp(-ct^{1/2}))^{-\nu}$	$c > 0, \nu > 0$

Table 2: Temporal building blocks proposed by [1]

Function	Parameters
$\psi_1(t) = (at^\alpha + 1)^\beta$	$a > 0, 0 < \alpha \leq 1, 0 \leq \beta \leq 1$
$\psi_2(t) = \ln(at^\alpha + b)/\ln(b)$	$a > 0, b > 1, 0 < \alpha \leq 1$
$\psi_3(t) = (at^\alpha + b)/(b(at^\alpha + 1))$	$a > 0, 0 < b \leq 1, 0 < \alpha \leq 1$

Table 3: Parameters of the covariance model given in equation (28) estimated by WLS.

$\sigma^2$	$\tau$	$c$	$\gamma$	$a$	$\alpha$	$\beta$
0.08	0.82	0.45	0.10	0.1	1.0	0.92

581 **List of Figures**

582     1	Raw data plots: (a) map of average log-concentrations of CHL 583       between 1998 and 2004, (b) map of average percentage of miss- 584       ing data for each location between 1998 and 2006, and (c) 585       time-series of the percentage of missing data over the Red Sea 586       between 1998 and 2004. . . . .	41
587     2	Log-CHL concentration time-series for every pixel of the do- 588       main (blue curves). The red curve plots the spatially averaged 589       log-concentrations and the black curve plots the spatial aver- 590       age of the seasonal component. Both are computed from the 591       data filled with DINEOF. . . . .	42
592     3	Anomalies estimated by DINEOF for the first 8 weeks of 1998. . . . .	43
593     4	Anomaly times series (panels b to d) at three locations as 594       indicated in panel (a) . . . . .	44
595     5	Description of the experimental setup. . . . .	45
596     6	RMS error over the cross-validation period for a varying num- 597       ber of smoothing iterations and number of EOFs. Crosses 598       indicate the minimal error for a given number of iterations. . . . .	46
599     7	EOF modes after the DINEOF data filling: (a) first spatial 600       EOF mode, (b) second spatial EOF mode, (c) first temporal 601       EOF mode, (d) second temporal EOF mode, (e) spatial av- 602       erage of the seasonal anomalies computed from the DINEOF 603       filled data, (f) cumulative sum of the percentage of variance 604       explained for the modes computed with DINEOF. . . . .	47
605     8	Fitting the covariance model to the space-time empirical covari- 606       ance matrix: (b) empirical space covariance function and 607       fitted space covariance functions, (c) empirical time covariance 608       function and fitted time covariance functions. . . . .	48
609     9	Contour plot of the empirical covariance function (dashed curves) 610       and of the fitted space-time covariance function (solid curves) . . . . .	49
611     10	Results of the Kalman filter: (a) time series of RMS errors 612       of forecasts for the covariance model compared with a pure 613       seasonal prediction, (b) distribution of the prediction RMS 614       error over the three-year cross-validation period. . . . .	50

615	11	Predictions for the period between 24 October 2006 and 31 October 2006 (a, b, c), and between 6 March 2006 and 14 March 2006 (d, e, f): (a, d) plot the seasonal prediction, (b, e) the model prediction with the Kalman filter, and (c, f) the observations. . . . .	51
620	12	Spatial averages over the validation period of: (a) the vari- ances of the prediction error resulting from by the reduced- order Kalman filter, and (b) the RMS errors. . . . .	52
625	13	Diagnostic statistics of the filter: (a) increment distribution in space averaged over the validation period, (b) cumulative aver- aged innovation over time, (c) distribution of the innovation values for the validation period. . . . .	53
630	14	Average correlation maps over the validation period between the observations and (a) the seasonal forecast, (b) the CHL forecasts (including the seasonal component), and (c) the filter analyses (including the seasonal component). . . . .	54

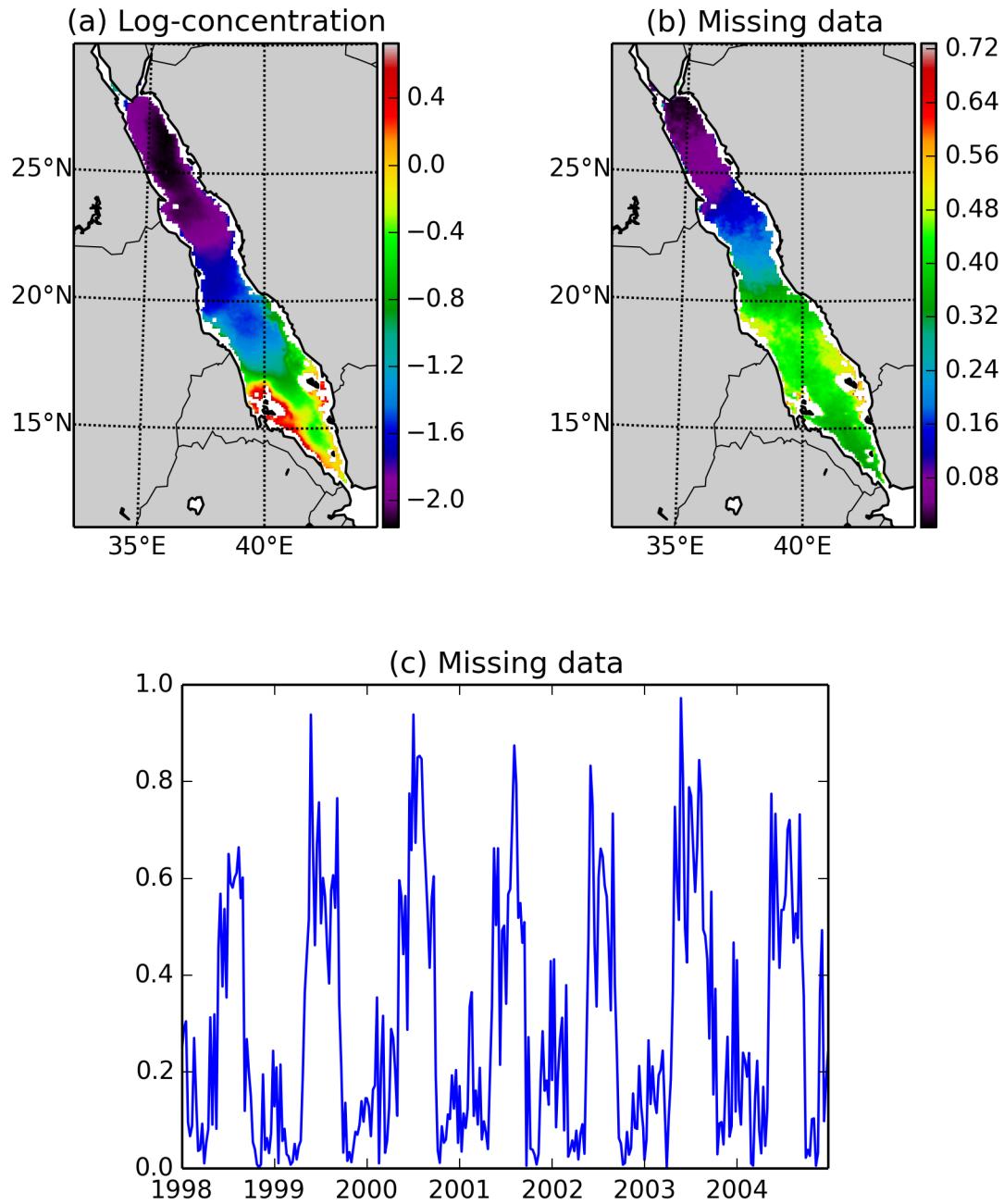


Figure 1: Raw data plots: (a) map of average log-concentrations of CHL between 1998 and 2004, (b) map of average percentage of missing data for each location between 1998 and 2006, and (c) time-series of the percentage of missing data over the Red Sea between 1998 and 2004.

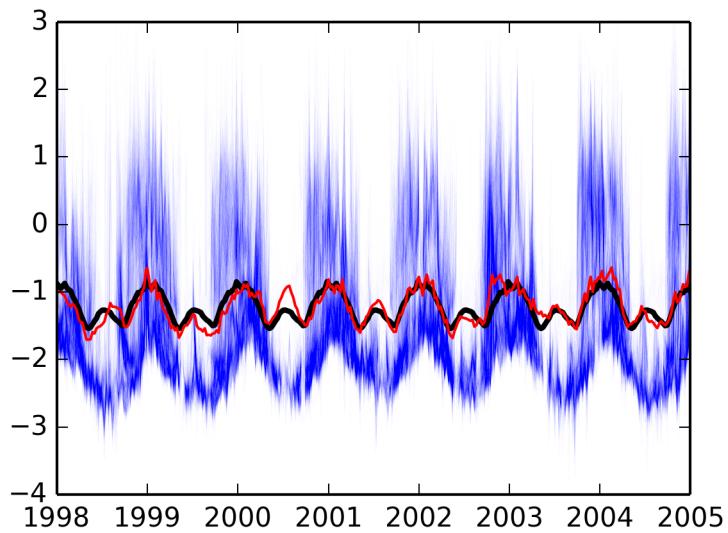


Figure 2: Log-CHL concentration time-series for every pixel of the domain (blue curves). The red curve plots the spatially averaged log-concentrations and the black curve plots the spatial average of the seasonal component. Both are computed from the data filled with DINEOF.

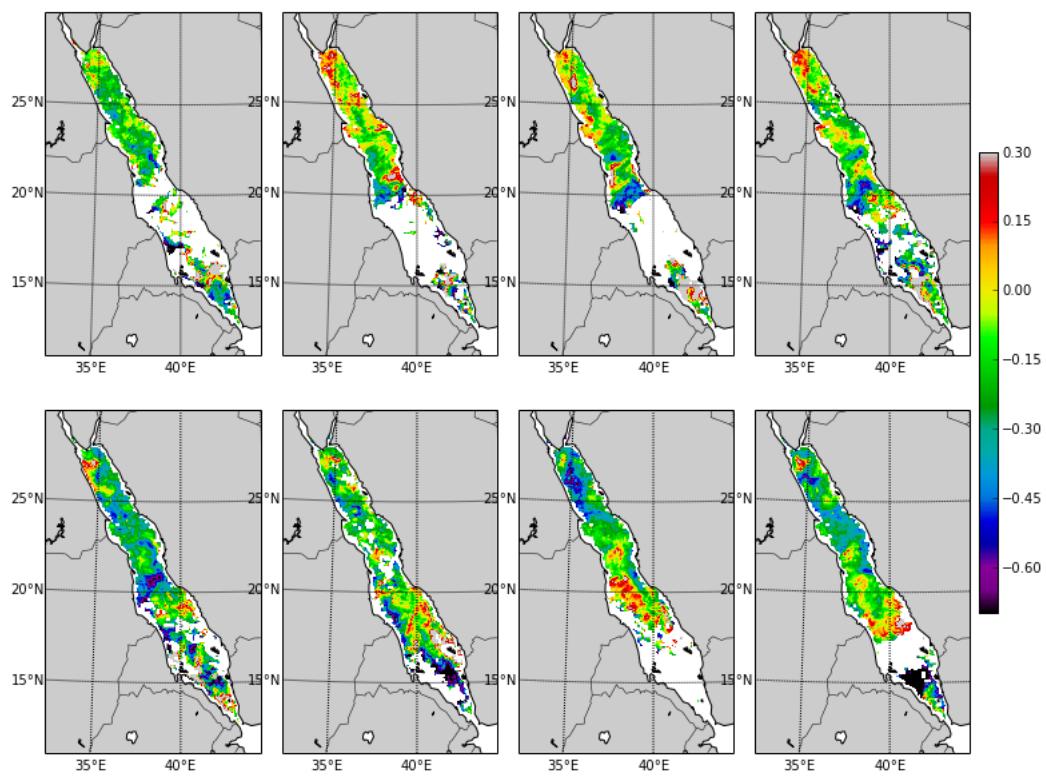


Figure 3: Anomalies estimated by DINEOF for the first 8 weeks of 1998.

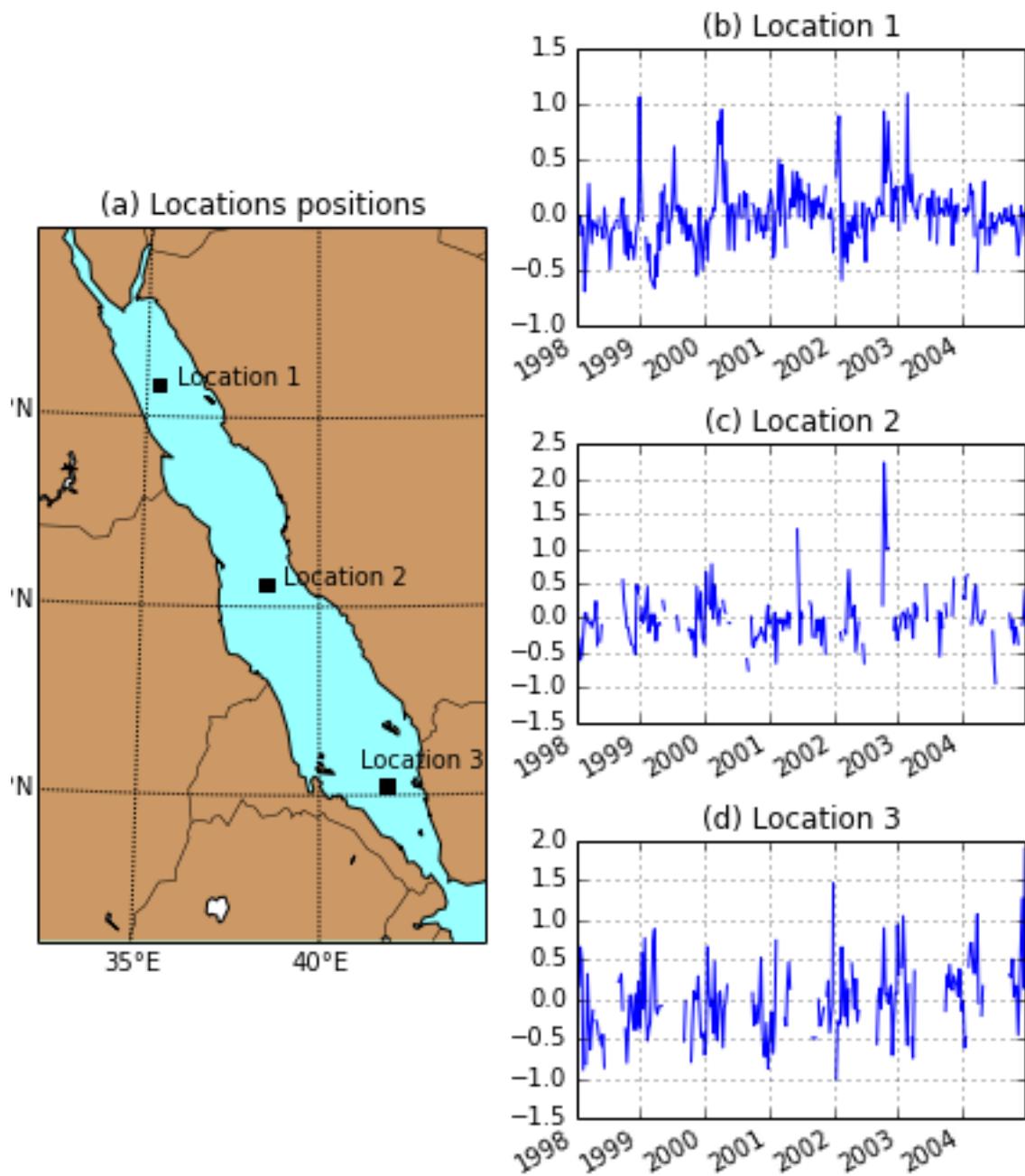


Figure 4: Anomaly times series (panels b to d) at three locations as indicated in panel (a)

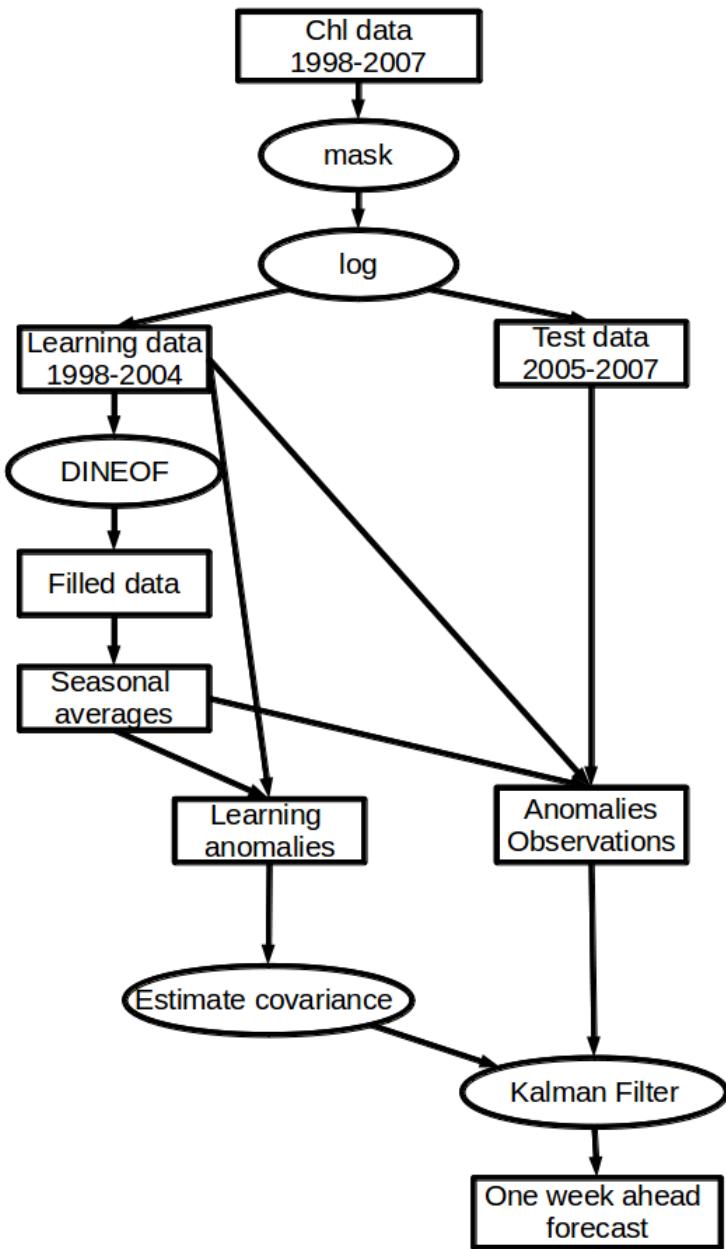


Figure 5: Description of the experimental setup.

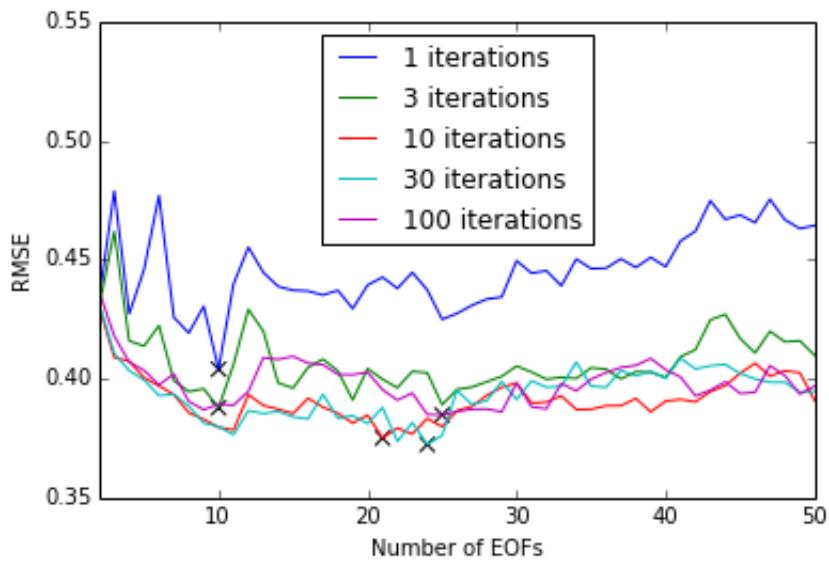


Figure 6: RMS error over the cross-validation period for a varying number of smoothing iterations and number of EOFs. Crosses indicate the minimal error for a given number of iterations.

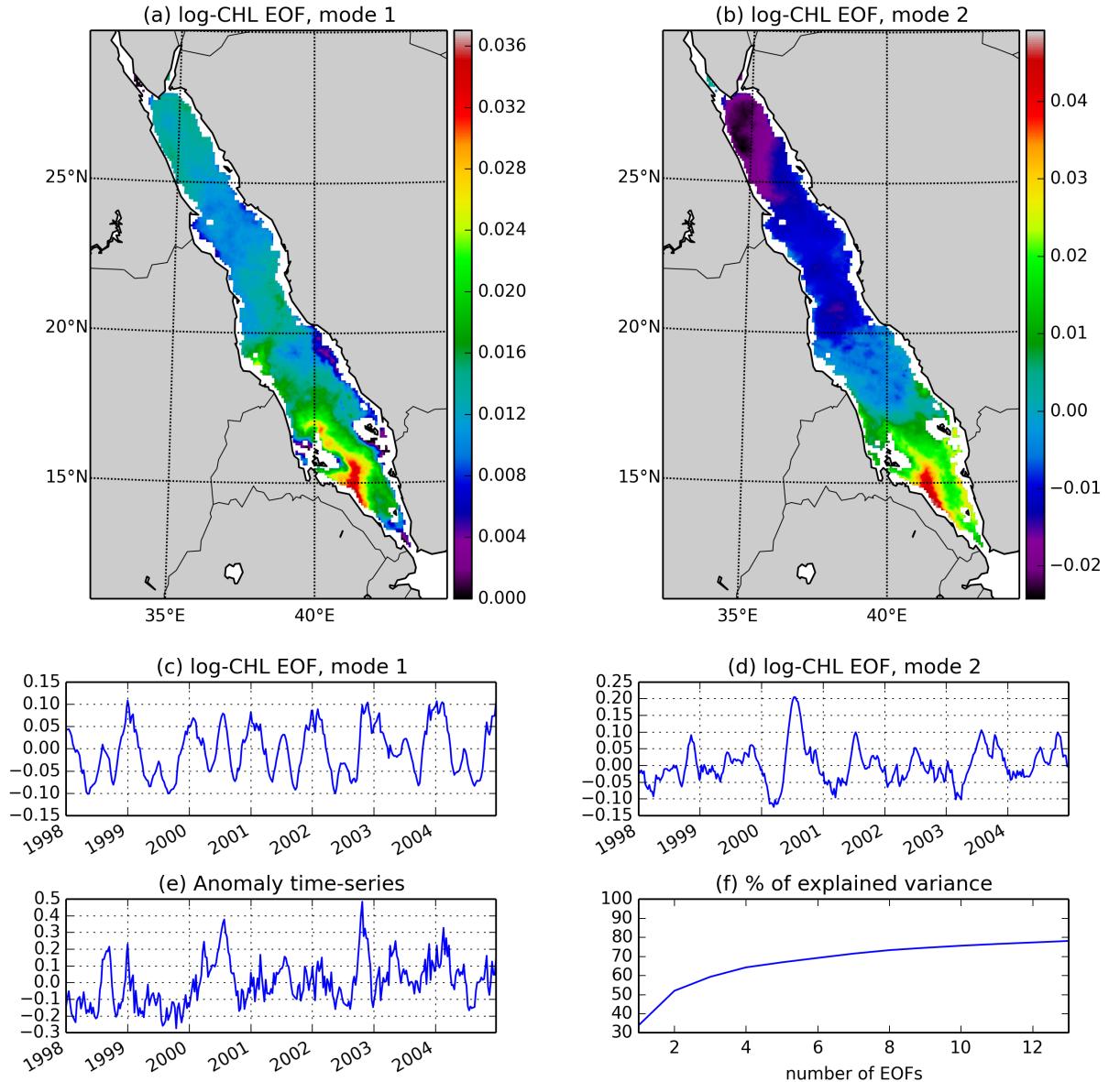


Figure 7: EOF modes after the DINEOF data filling: (a) first spatial EOF mode, (b) second spatial EOF mode, (c) first temporal EOF mode, (d) second temporal EOF mode, (e) spatial average of the seasonal anomalies computed from the DINEOF filled data, (f) cumulative sum of the percentage of variance explained for the modes computed with DINEOF.

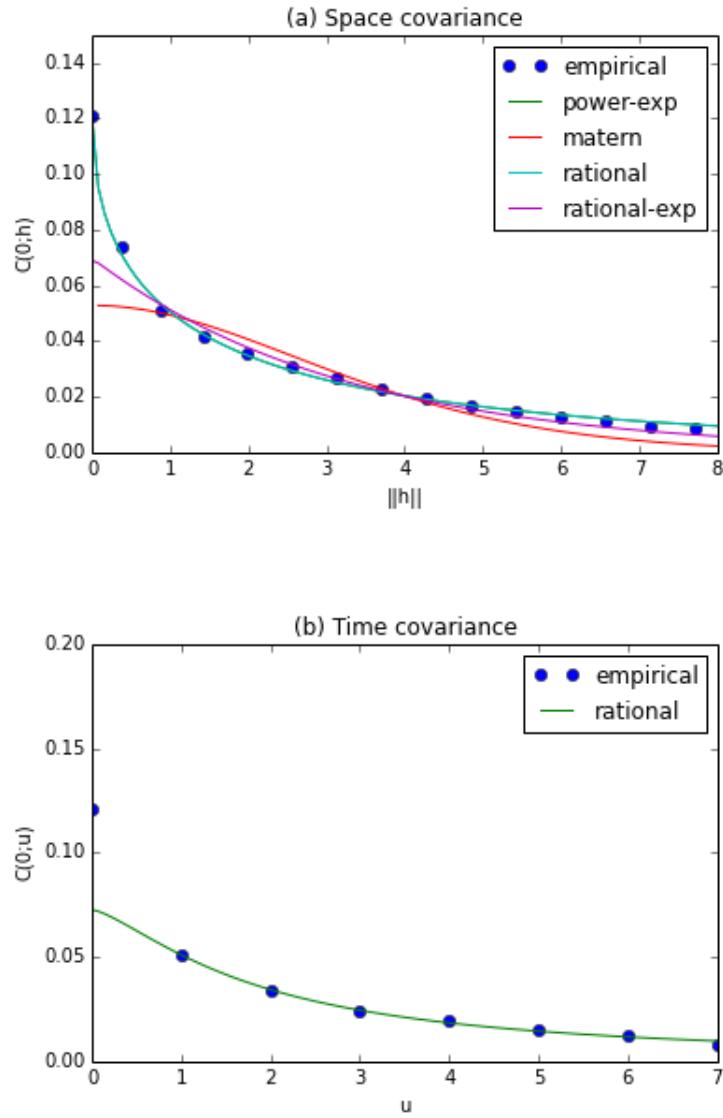


Figure 8: Fitting the covariance model to the space-time empirical covariance matrix: (b) empirical space covariance function and fitted space covariance functions, (c) empirical time covariance function and fitted time covariance functions.

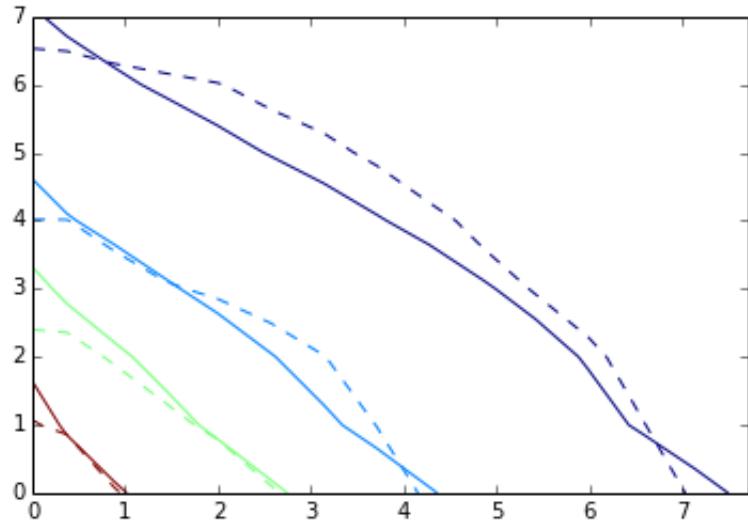


Figure 9: Contour plot of the empirical covariance function (dashed curves) and of the fitted space-time covariance function (solid curves)

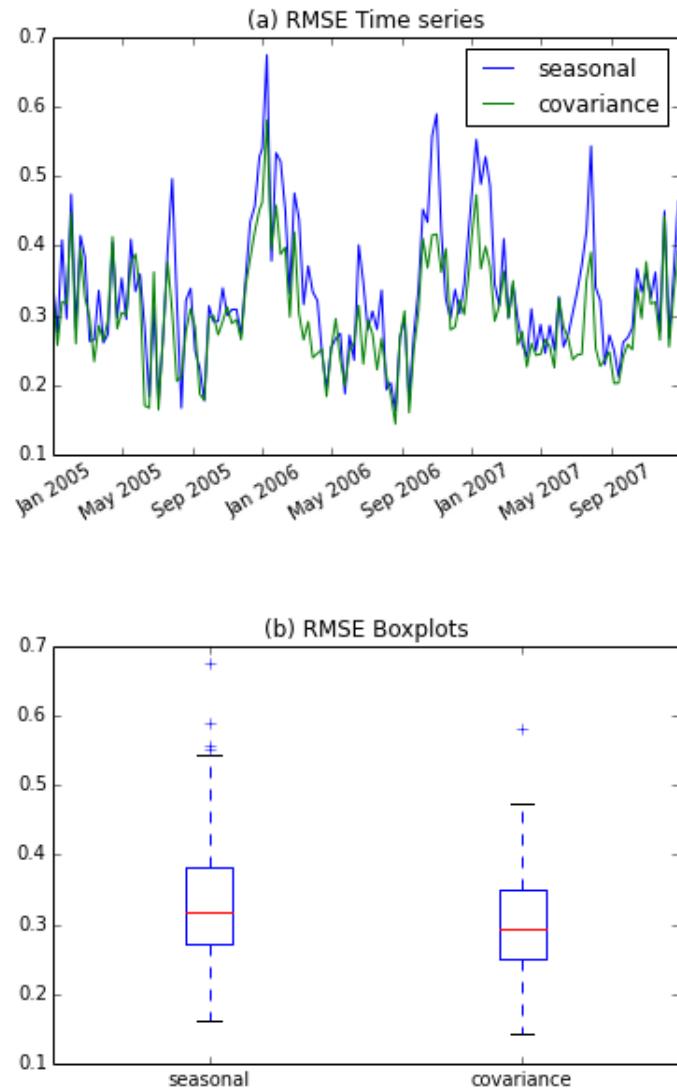


Figure 10: Results of the Kalman filter: (a) time series of RMS errors of forecasts for the covariance model compared with a pure seasonal prediction, (b) distribution of the prediction RMS error over the three-year cross-validation period.

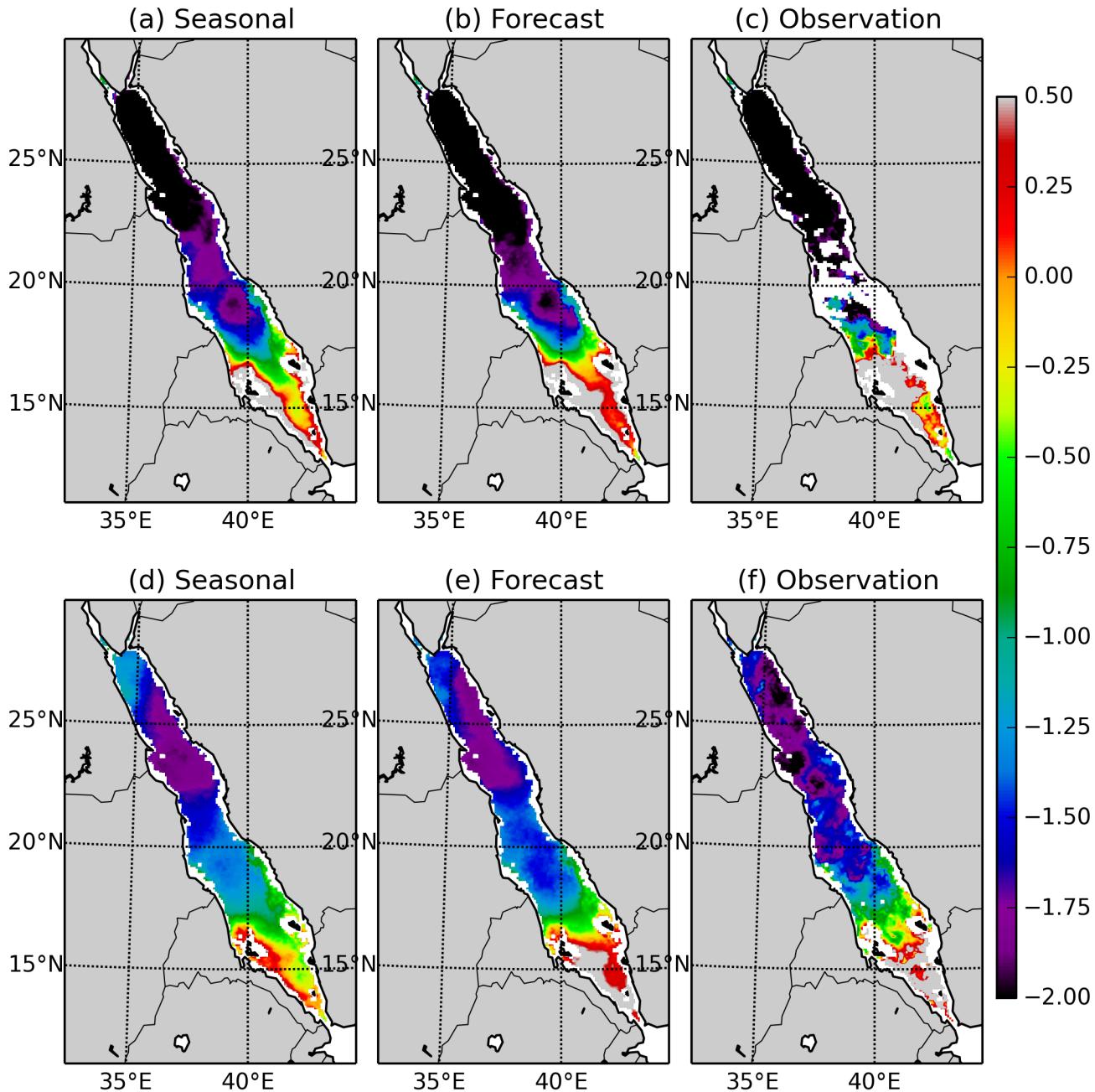


Figure 11: Predictions for the period between 24 October 2006 and 31 October 2006 (a, b, c), and between 6 March 2006 and 14 March 2006 (d, e, f): (a, d) plot the seasonal prediction, (b, e) the model prediction with the Kalman filter, and (c, f) the observations.

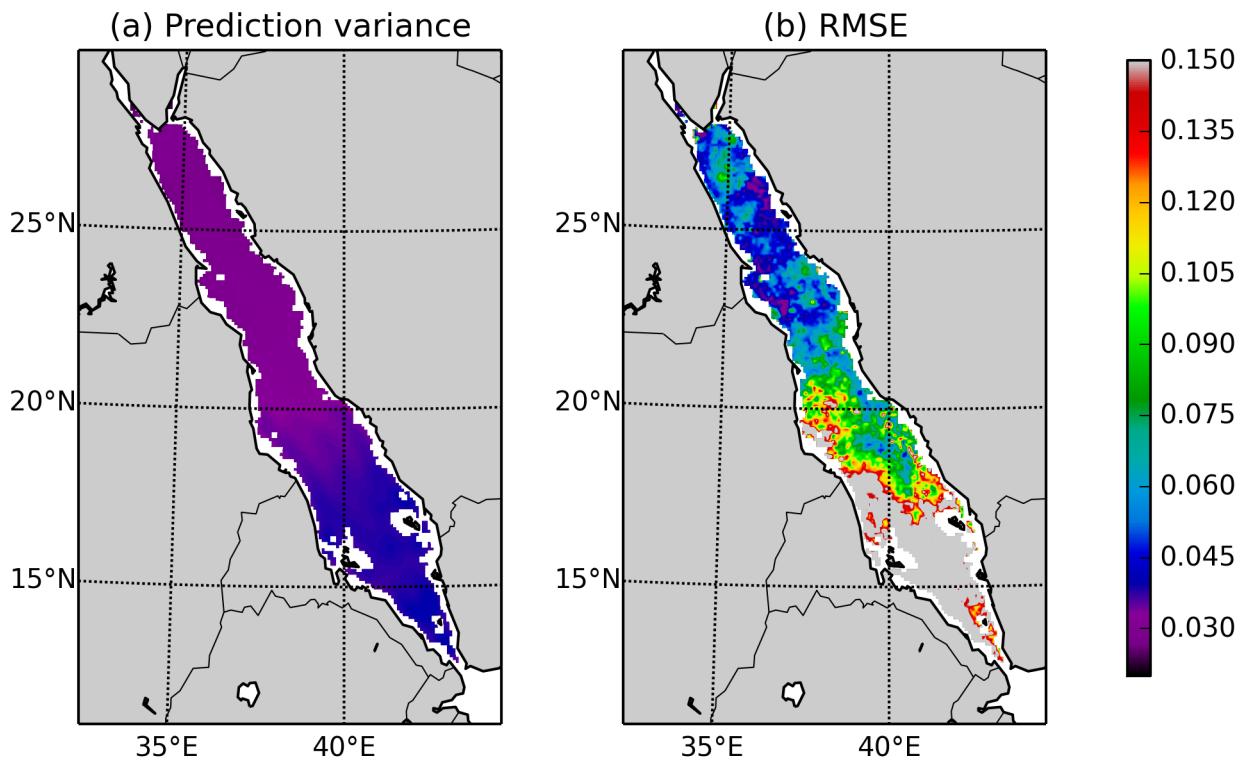


Figure 12: Spatial averages over the validation period of: (a) the variances of the prediction error resulting from by the reduced-order Kalman filter, and (b) the RMS errors.

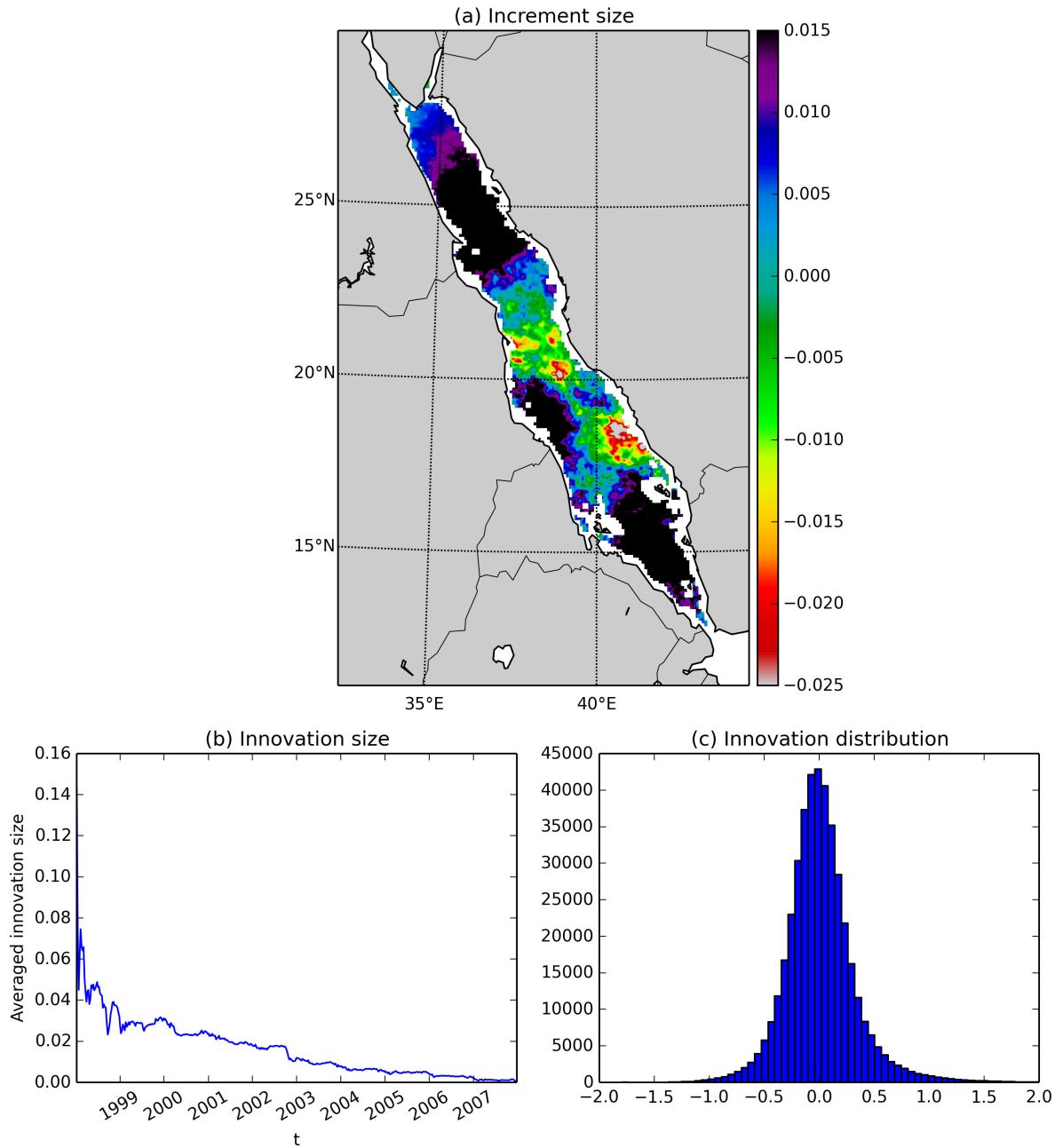


Figure 13: Diagnostic statistics of the filter: (a) increment distribution in space averaged over the validation period, (b) cumulative averaged innovation over time, (c) distribution of the innovation values for the validation period.

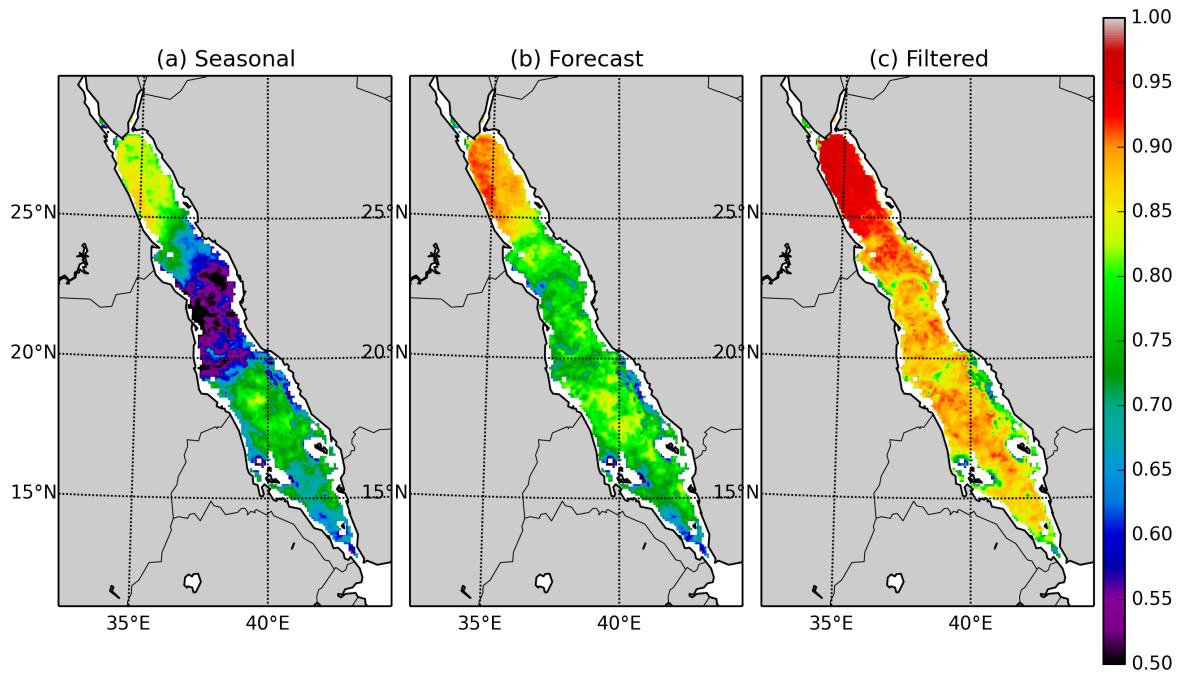


Figure 14: Average correlation maps over the validation period between the observations and (a) the seasonal forecast, (b) the CHL forecasts (including the seasonal component), and (c) the filter analyses (including the seasonal component).