

# Data assimilation in a coupled physical-biogeochemical model of the California Current System using an incremental lognormal 4-dimensional variational approach: Part 2, Joint physical and biological data assimilation twin experiments

Hajoon Song<sup>a,\*</sup>, Christopher A. Edwards<sup>b</sup>, Andrew M. Moore<sup>b</sup>, Jerome Fiechter<sup>b</sup>

<sup>a</sup>*Department of Earth, Atmospheric and Planetary Sciences, Massachusetts Institute of Technology, 77 Massachusetts Avenue, Cambridge, MA 02137, U.S.A.*

<sup>b</sup>*Ocean Sciences Department, University of California, 1156 High Street, Santa Cruz, CA 96064, U.S.A.*

---

## Abstract

Coupled physical and biological data assimilation is performed within the California Current System using model twin experiments. The initial condition of physical and biological variables is estimated using the four-dimensional variational (4DVar) method under the Gaussian and lognormal error distributions assumption, respectively. Errors are assumed to be independent, yet variables are coupled by assimilation through model dynamics. Using a nutrient-phytoplankton-zooplankton-detritus (NPZD) model coupled to an ocean circulation model (the Regional Ocean Modeling System, ROMS), the coupled data assimilation procedure is evaluated by comparing results to experiments with no assimilation and with assimilation of physical data and

---

\*Corresponding author, Tel. : +1 617 253 0098  
*Email address:* [hajsong@mit.edu](mailto:hajsong@mit.edu) (Hajoon Song)

biological data separately. Independent assimilation of physical (biological) data reduces the root-mean-squared error (RMSE) of physical (biological) state variables by more than 56% (43%). However, the improvement in biological (physical) state variables is less than 7% (13%). In contrast, coupled data assimilation improves both physical and biological components by 57% and 49%, respectively. Coupled data assimilation shows robust performance with varied observational errors, resulting in significantly smaller RMSEs compared to the free run. It still produces the estimation of observed variables better than that from the free run even with the physical and biological model error, but leads to higher RMSEs for unobserved variables. A series of twin experiments illustrates that coupled physical and biological 4DVar assimilation is computationally efficient and practical, capable of providing the reliable estimation of the coupled system with the same and ready to be examined in a realistic configuration.

*Keywords:* Coupled data assimilation, Biogeochemical model, 4DVAR, California Current System Keyword

---

## 1. Introduction

Marine ecosystem and biogeochemical models coupled to realistic ocean circulation models are applied routinely today for an extensive range of studies, such as primary production (Franks and Chen, 2001), ecosystem phenology (Chenillat et al., 2013), biogeography (Follows et al., 2007, Goebel et al., 2010), nutrient cycling (Fennel, 2010), air-sea carbon exchange (Chai et al., 2009) and climate change (Cox et al., 2000, Behrenfeld et al., 2006, Stock et al., 2011). Despite real advances in the representation of complex biolog-

9 ical interactions and improvements in physical circulation modeling, many  
10 factors still contribute to errors in model output, including imperfect param-  
11 eterization of biological and physical processes at both resolved and subgrid  
12 scales. In realistic applications in which ocean state estimates at particular  
13 times are sought, for example as part of an ocean observing system, addi-  
14 tional errors result from uncertainties in ocean model initial conditions and  
15 applied forcing.

16 One approach to improving model fidelity for ocean state estimation is  
17 through data assimilation in which model estimates are constrained through  
18 adjustment of control variables to better match available observations. Devel-  
19 opments in data assimilation in physical oceanography for over two decades  
20 now provide many estimates of the physical ocean state at global (Behringer  
21 et al., 1998, Bell et al., 2000, Chassignet et al., 2007, Köhl et al., 2007,  
22 Balmaseda et al., 2008, Carton and Giese, 2008) and regional (Oke et al.,  
23 2008, Cummings et al., 2009, Broquet et al., 2009, Shulman et al., 2009, Ku-  
24 rapov et al., 2011, Matthews et al., 2012, Sakov et al., 2012) scales. Data  
25 assimilation techniques have been developed to a lesser extent in biological  
26 oceanography, though their application has been used for both the determi-  
27 nation of poorly known model parameters (e.g., Matear, 1995, Spitz et al.,  
28 1998, Losa et al., 2004, Mattern et al., 2012, Roy et al., 2012, Doron et al.,  
29 2013, Simon et al., 2015) as well as quantitative improvement in modeled  
30 biological fields (see Gregg (2008) and Edwards et al. (2015) for reviews).

31 Data assimilation in physical-biological coupled systems has focused gener-  
32 ally on either physical or biological data assimilation in isolation. Better  
33 representation of the ocean circulation by physical data assimilation is ex-

34 pected to improve the distribution of biological variables. Studies focusing on  
35 the impact of physical data assimilation on biological fields include Oschlies  
36 and Garçon (1998), Miller et al. (2000), Berline et al. (2007) and Fiechter  
37 et al. (2011). For example, Oschlies and Garçon (1998) assimilate satellite  
38 estimates of sea surface height (SSH) over the North Atlantic to improve  
39 the eddy representation in the physical model, which provides currents for a  
40 coupled biological model. They report that the nitrate flux into the euphotic  
41 zone is increased by improving the underestimated mesoscale eddy activity in  
42 the free simulation. Raghukumar et al. (2015) present a counter example in  
43 which physical data assimilation alone can drive spurious nutrient fluxes into  
44 the euphotic zone degrading ecosystem model performance. Improvements in  
45 biological fields have also resulted from assimilation of biological fields alone,  
46 where physical fields have been assumed a priori to be sufficiently accurate or  
47 already modified through physical data assimilation (e.g., Friedrichs, 2001,  
48 Garcia-Gorriz et al., 2003, Natvik and Evensen, 2003, Hoteit et al., 2005,  
49 Triantafyllou et al., 2007, Ciavatta et al., 2011, Rousseaux and Gregg, 2012,  
50 Ford et al., 2012, Hu et al., 2012). For example, Garcia-Gorriz et al. (2003)  
51 assimilate the Sea-viewing Wide Field-of-view Sensor (SeaWiFS) data into  
52 an NPZ model coupled to a 3D ocean model over the Adriatic Sea, improving  
53 ecosystem parameters to reduce misfits between observations and the model  
54 output. Natvik and Evensen (2003) use SeaWiFS data to fit an 11-component  
55 biochemical model coupled to a 3D ocean circulation model configured over  
56 the North Atlantic. Using an ensemble Kalman filter approach, they adjust  
57 state variables and show that multivariate biochemical data assimilation can  
58 not only improve the representation of an observed variable but also reduce

59 the variance of unobserved variables.

60 A few studies have addressed the assimilation of both physical and bi-  
61 ological data into coupled models and its advantages over the fit to either  
62 physical or biological data alone. Using an optimal interpolation approach  
63 applied to the Gulf Stream region, Anderson et al. (2000) find that dynam-  
64 ically consistent physical and biological fields created through joint physi-  
65 cal and biological assimilation are superior to those obtained through either  
66 physical or biological data assimilation alone. In a sequential data assimila-  
67 tion study of the North Atlantic, Ourmières et al. (2009) find that ecosystem  
68 state estimates are improved through assimilation of both physical (sea sur-  
69 face temperature (SST), SSH and climatological temperature (T) and salin-  
70 ity (S)) and biological (nitrate climatology) data more than through physical  
71 data assimilation alone. Indeed in their results, physical data assimilation  
72 alone may degrade ecosystem estimates, depending on the accuracy of the  
73 nitrate background state. In a study of Monterey Bay, Shulman et al. (2013)  
74 report that substantial improvement in biological estimates did not result  
75 from physical assimilation alone but required assimilation of biological fields  
76 or updates to biological fields (nitrate) through statistical relations. Simon  
77 et al. (2015) perform two steps of assimilation using an ensemble Kalman  
78 filter in coupled physical and biological state estimation. After fitting the  
79 model to SST, along-track sea level anomalies and ice concentration obser-  
80 vations, they assimilate 8-day composite chlorophyll data to estimate the  
81 biological states and parameters in log-transformed space. A stronger error  
82 reduction results from assimilating all observations, but physical and biolog-  
83 ical components are independent during assimilation.

84 To date biological assimilation efforts on adjusting state variables have  
85 largely used sequential methods, based on optimal interpolation or Kalman  
86 filter approaches. In particular, Bertino et al. (2003) applied Gaussian anamor-  
87 phosis to biogeochemical variables with the non-Gaussian distributions in  
88 the ensemble Kalman filter. This transformation satisfies the assumption of  
89 Gaussian error distribution and they show promising results in fitting 1-D  
90 numerical ecosystem model to observations. Recently, a four-dimensional  
91 variational assimilation method appropriate for ocean ecosystem variables  
92 was studied in an idealized 1-dimensional context (Song et al., 2012). This  
93 method accounts for the non-Gaussian statistics of ecosystem variables by as-  
94 suming lognormal statistics following Fletcher and Zupanski (2006). In Song  
95 et al. (2016a), this approach is modified and implemented within a realis-  
96 tic ocean circulation model (ROMS; the Regional Ocean Modeling System).  
97 The modification includes a linearization of the log-transformation function  
98 to enable efficient searching for the cost function minimum. Although the lin-  
99 earization requires the exclusion of observations whose values substantially  
100 exceed the background state, the modified log-transformed 4DVar outper-  
101 forms the conventional 4DVar (which assumes Gaussian error distributions)  
102 both in terms of RMSE and non-negativity in a series of model twin experi-  
103 ments configured for the California Current System (Song et al., 2016a).

104 In this study, we extend that work by developing the ability to jointly  
105 assimilate physical and biological data within ROMS. The four-dimensional  
106 variational (4DVar) method provides dynamically consistent state estimates  
107 within each assimilation cycle. Coupled dynamics within the tangent linear  
108 and adjoint models of the 4DVar system have the potential to enable both

109 physical data to improve biological fields and biological data to improve phys-  
110 ical fields. In most sequential data assimilation methods the latter is possible  
111 only through statistical adjustments and not directly through model dynam-  
112 ics. We investigate the advantage of coupled data assimilation using both  
113 physical and biological data by comparing results from multiple runs: no  
114 data assimilation, physical data assimilation alone, biological data assimila-  
115 tion alone, and joint physical and biological assimilation. Using model twin  
116 experiments, we fit the coupled model to pseudo observations of SST, SSH,  
117 in situ temperature and salinity, and surface chlorophyll data. The statisti-  
118 cal analyses highlight the advantage of the coupled data assimilation in the  
119 present model configuration.

120 The organization of this paper is as follows. A brief introduction to phys-  
121 ical and biological variational data assimilation method is given in section  
122 2. Section 3 describes the twin experiment design to evaluate the coupled  
123 data assimilation system and its performance is presented in section 4. We  
124 summarize results and provide a discussion in section 5 to end the paper.

## 125 **2. Coupled variational data assimilation**

126 It is desirable to perform physical and biological data assimilation simul-  
127 taneously rather than independently because (a) it is computationally more  
128 efficient to carry out a single assimilation procedure with a larger model  
129 than to perform two sequential but smaller assimilation operations and (b)  
130 the physical and biological fields are coupled through model dynamics. This  
131 coupling enables biological observations to constrain physical fields and vice  
132 versa, in principle leading to improved state estimates over results from in-

133 dependent (i.e., uncoupled) physical and biological assimilation. In practice,  
134 however, the beneficial results of coupled assimilation must be demonstrated,  
135 and issues such as model error and observational uncertainty may limit the  
136 ultimate improvements obtained through coupled state estimation.

137 Fundamentally, the basic error assumptions for physical and biological  
138 data assimilation are different. Errors in physical variables are assumed to  
139 be Gaussian distributed, whereas errors in biological variables are better  
140 represented by lognormal distributions (Campbell, 1995, Simon and Bertino,  
141 2009); as a result, the assumption of Gaussian-distributed errors is appro-  
142 priate for physical but not biological variables. Here, we present a method  
143 for combining the two different 4DVar approaches by following Fletcher and  
144 Zupanski (2006), Fletcher (2010) and Fletcher and Jones (2014). Although  
145 the state vector could include initial conditions, boundary conditions, sur-  
146 face forcing fields, and biological parameters, we consider here for simplicity  
147 a control vector consisting only of the initial state  $\mathbf{x}_0$ .

148 We first adapt the incremental form of 4DVar following Song et al. (2016a).  
149 If the posterior initial condition  $\mathbf{x}_0$  is written as the sum of the background  
150 state  $\mathbf{x}_{b,0}$  and a (small) increment  $\delta\mathbf{x}_0$ , then  $\mathbf{x}_i^o = \mathcal{H}_i(\mathcal{M}_{i,0}(\mathbf{x}_{b,0} + \delta\mathbf{x}_0)) \approx$   
151  $\mathcal{H}_i(\mathcal{M}_{i,0}(\mathbf{x}_{b,0})) + \mathbf{H}_i\mathbf{M}_{i,0}\delta\mathbf{x}_0$ , where the nonlinear model  $\mathcal{M}_{i,0}$  integrates the  
152 state vector from  $t = t_0$  to  $t = t_i$ , and the observation operator  $\mathcal{H}_i$  maps  
153 model states to observation space and  $\mathbf{x}_{b,0}$  is the background state vector.  
154 The matrices  $\mathbf{H}_i$  and  $\mathbf{M}_{i,0}$  are the tangent linear forms of  $\mathcal{H}_i$  and  $\mathcal{M}_{i,0}$ , respec-  
155 tively. In this case, the cost function  $J_G$  appropriate for Gaussian-distributed



156 variables becomes

$$\begin{aligned}
J_G(\delta \mathbf{x}_0) &= \frac{1}{2} \delta \mathbf{x}_0^T \mathbf{B}_G^{-1} \delta \mathbf{x}_0 \\
&+ \frac{1}{2} \sum_{i=1}^{N_o} (\mathbf{d}_{g,i} - \mathbf{H}_i \mathbf{M}_{i,0} \delta \mathbf{x}_0)^T \mathbf{R}_{G,i}^{-1} (\mathbf{d}_{g,i} - \mathbf{H}_i \mathbf{M}_{i,0} \delta \mathbf{x}_0), \quad (1)
\end{aligned}$$

157 where  $\mathbf{d}_{g,i} = \mathbf{y}_i - \mathbf{x}_{b,i}^o = \mathbf{y}_i - \mathcal{H}_i(\mathcal{M}_{i,0}(\mathbf{x}_{b,0}))$  define the innovations.

158 Similarly, the cost function  $J_L$  for lognormally-distributed variables is  
159 expressed in terms of increments  $\delta \mathbf{g}_0 = \ln \mathbf{x}_0 - \ln \mathbf{x}_{b,0}$  as

$$\begin{aligned}
J_L(\delta \mathbf{g}_0) &= \frac{1}{2} \delta \mathbf{g}_0^T \mathbf{B}_L^{-1} \delta \mathbf{g}_0 \\
&+ \frac{1}{2} \sum_{i=1}^{N_o} (\mathbf{d}_{l,i} - \mathbf{O}_{L,i} \mathbf{H}_i \mathbf{M}_{i,0} \mathbf{X}_L \delta \mathbf{g}_0)^T \mathbf{R}_{L,i}^{-1} \\
&\quad (\mathbf{d}_{l,i} - \mathbf{O}_{L,i} \mathbf{H}_i \mathbf{M}_{i,0} \mathbf{X}_L \delta \mathbf{g}_0), \quad (2)
\end{aligned}$$

160 where  $\mathbf{d}_{l,i} = \ln \mathbf{y}_i - \ln \mathbf{x}_{b,i}^o = \ln \mathbf{y}_i - \ln(\mathcal{H}_i(\mathcal{M}_{i,0}(\mathbf{x}_{b,0})))$ . The diagonal matrices  
161  $\mathbf{O}_{L,i}$  and  $\mathbf{X}_L$  are introduced during the linearization of ln and exp function  
162 (Song et al., 2016a) where specifically,  $\mathbf{O}_{L,i} = \text{diag}[(\mathbf{x}_{b,i}^o)_1, (\mathbf{x}_{b,i}^o)_2, \dots, (\mathbf{x}_{b,i}^o)_{m_i}]^{-1}$   
163 and  $\mathbf{X}_L = \text{diag}[(\mathbf{x}_{b,0})_1, (\mathbf{x}_{b,0})_2, \dots, (\mathbf{x}_{b,0})_{n_l}]$ . A total of  $m_i$  lognormally-  
164 distributed observations exist at  $t = t_i$ , and a total of  $n_l$  lognormal variables  
165 exist within the model.

166 The state vector increment in the coupled, physical and biological system  
167 is defined as  $\delta \mathbf{z} = [\delta \mathbf{x}_0^T \quad \delta \mathbf{g}_0^T]^T$ . The state vectors for physical and biological  
168 variables have dimensions  $(n_g \times 1)$  and  $(n_l \times 1)$ , respectively. Hence the size  
169 of  $\delta \mathbf{z}_0$  is simply  $(n \times 1)$ , where  $n = n_g + n_l$ . Then a compact form of the cost

170 function can be written as

$$\begin{aligned}
J(\delta \mathbf{z}_0) &= \frac{1}{2} \delta \mathbf{z}_0^T \mathbf{B}^{-1} \delta \mathbf{z}_0 \\
&+ \frac{1}{2} \sum_{i=1}^{N_o} (\mathbf{d}_i - \mathbf{O}_i^{-1} \mathbf{H}_i \mathbf{M}_{i,0} \mathbf{X} \delta \mathbf{z}_0)^T \mathbf{R}_i^{-1} \\
&(\mathbf{d}_i - \mathbf{O}_i^{-1} \mathbf{H}_i \mathbf{M}_{i,0} \mathbf{X} \delta \mathbf{z}_0), \tag{3}
\end{aligned}$$

171 where  $\mathbf{d}_i^T = [\mathbf{d}_{g,i}^T \ \mathbf{d}_{b,i}^T]$ ,  $\mathbf{O}_i = \text{diag}[1, 1, \dots, 1, (\mathbf{x}_{b,i}^o)_1, (\mathbf{x}_{b,i}^o)_2, \dots, (\mathbf{x}_{b,i}^o)_{m_i}]$   
172 and  $\mathbf{X} = \text{diag}[1, 1, \dots, 1, (\mathbf{x}_{b,0})_1, (\mathbf{x}_{b,0})_2, \dots, (\mathbf{x}_{b,0})_{n_i}]$ . Error covariances  
173  $\mathbf{B}$  and  $\mathbf{R}_i$  consist of error covariances for physical and biological components  
174 and their cross covariances.

175 The gradient of  $J(\delta \mathbf{z}_0)$  with respect to  $\delta \mathbf{z}_0$  is given by

$$\frac{\partial J}{\partial \delta \mathbf{z}_0} = \mathbf{B}^{-1} \delta \mathbf{z}_0 - \mathbf{X}^T \sum_{i=1}^{N_o} \mathbf{M}_{0,i}^T \mathbf{H}_i^T \mathbf{O}_i^{-T} \mathbf{R}_i^{-1} (\mathbf{d}_i - \mathbf{O}_i^{-1} \mathbf{H}_i \mathbf{M}_{i,0} \mathbf{X} \delta \mathbf{z}_0), \tag{4}$$

176 and we seek a solution  $\delta \mathbf{z}_0$  that satisfies  $\partial J / \partial \delta \mathbf{z}_0 = 0$ . The optimal  $\delta \mathbf{z}_0$   
177 is identified iteratively by applying a conjugate gradient descent algorithm  
178 using the Lanczos formulation (Moore et al., 2011c).

### 179 **3. Experiment design for the coupled data assimilation system** 180 **evaluation**

181 We evaluate the performance of the new system by comparing results  
182 from multiple data assimilation experiments with a free run. For clarity of  
183 description, we distinguish between the coupled nonlinear model and the  
184 coupled data assimilation system by referring to the former as the forward  
185 model with no data assimilation. The free run results exclusively from the  
186 integration of the forward model.

187 We perform three data assimilation runs: the physical data assimilation  
188 (PDA) run, the biological data assimilation (BDA) run and the coupled  
189 data assimilation (PBDA) run. In the PDA run, physical data are used to  
190 constrain both physical and biological variables in the forward model. In the  
191 BDA run, only biological data are used to constrain physical and biological  
192 variables in the forward model. Third, both physical and biological data are  
193 used to constrain the physical and biological variables in the coupled model,  
194 and is referred to as the PBDA run. The PDA (BDA) applies Gaussian  
195 (Lognormal) 4DVar to fit the coupled model to the data. The PBDA fits  
196 both physical and biological data into the coupled model using a hybrid  
197 Gaussian and lognormal 4DVar approach.

198 Some additional experiments were performed to better evaluate the im-  
199 pact of model dynamics within the coupled assimilative system. Specifically,  
200 we modify the BDA experiment by reducing the control vector to include  
201 only biological variables alone (BDAb) or physical variables only (BDAp).  
202 Similarly, we consider the PDA experiment, but with adjustments to only  
203 the physical (PDAp) or biological (PDAb) variables. While physical state  
204 variables in the forward model clearly influence biological variables (e.g.,  
205 through transport and mixing), biological variables generally do not alter  
206 physical variables in the forward model. In nature, chlorophyll pigments  
207 quantitatively impact light absorption and thus heat flux within the water  
208 column (Morel, 1988, Lewis et al., 1990, Frouin and Iacobellis, 2002, Mur-  
209 tugudde et al., 2002, Park et al., 2015), but this feedback is not included in  
210 the present model implementation. As a result, any misfit with respect to  
211 physical data cannot be reduced by adjusting the initial conditions of biolog-

212 ical variables. Therefore, the physical data misfit is identical for PDAp and  
213 PDA, and for PDAb and the free run.

### 214 *3.1. Model*

215 We use a Nutrient-Phytoplankton-Zooplankton-Detritus (NPZD) biolog-  
216 ical model coupled to ROMS, a 3-dimensional ocean circulation model, con-  
217 figured for the California Current System. This implementation has been  
218 applied repeatedly as a useful testbed for various developments of the ROMS  
219 4DVar system (Broquet et al., 2009, 2011, Moore et al., 2011b,a). The model  
220 domain extends from the middle of the Baja Peninsula to the Washington  
221 coast and offshore to 137W. The horizontal model resolution is 1/30 degree,  
222 and it includes 30 terrain-following levels in the vertical. The configuration  
223 used here is identical to that testing the lognormal 4DVar in isolation and  
224 described in Song et al. (2016a), which provides additional details of the  
225 configuration, including the parameters used for the NPZD model.

226 Model twin experiments are an excellent way to evaluate the performance  
227 of data assimilation schemes because the true state is known exactly and the  
228 error statistics can be controlled. A 4-year forward simulation, begun on  
229 January 1<sup>st</sup>, 2001, represents the “true” ocean state. A data assimilated run  
230 described by Broquet et al. (2009) provides the physical initial conditions.  
231 Biological initial conditions were obtained from the final state of a 45-year  
232 forward spin-up run described in Song et al. (2016a). Surface forcing and  
233 boundary conditions were derived from the output of the Coupled Ocean  
234 Atmosphere Mesoscale Prediction System (COAMPS) (Doyle et al., 2009)  
235 and the Simple Ocean Data Assimilation (SODA) (Carton and Giese, 2008)  
236 product, respectively.

237 Although real data is not used in the present experiments, we present an  
238 evaluation of the forward simulation using satellite derived estimates of sea  
239 surface temperature (AVHRR Pathfinder V5 SST, 0.44 degree resolution)  
240 and surface chlorophyll-a (SeaWiFS, 0.036 degree resolution) obtained from  
241 <http://las.pfeg.noaa.gov/oceanWatch>. Model chlorophyll is estimated using  
242 a constant carbon to chlorophyll ratio of 50 g C (g chl)<sup>-1</sup> and a Redfield  
243 ratio to convert model units of nitrogen to carbon. The model yields phyto-  
244 plankton bloom-like patterns, intensity and spatial distribution comparable  
245 to satellite data (Song et al., 2016a). Monthly average fields are used to cal-  
246 culate the bias (model minus data), normalized by the standard deviation,  
247 and correlation coefficient (Figure 1). The standard deviation is estimated  
248 at each grid cell using the output from the model spin-up. Although SST in  
249 the simulation has a warm bias overall (1 °C), the correlation coefficient ( $r$ )  
250 is very high ( $\bar{r} = 0.92$ ), indicating a good representation of the annual cycle  
251 in the model. Surface chlorophyll-a is biased low offshore and very near the  
252 coast north of 44°N, and biased high along the northern and central Cali-  
253 fornia coast out into the coastal transition zone (Brink and Cowles, 1991).  
254 On average, the model is biased low by approximately 0.5 mg m<sup>-3</sup>. The  
255 correlation coefficient for chlorophyll-a reveals generally positive values over  
256 the whole domain.

### 257 3.2. Data

258 Physical and biological data are sampled from the true state. SSH and  
259 SST are observed at all grid points once per day (we assume no data dropouts  
260 due to cloud cover). In situ temperature and salinity profiles are obtained  
261 at times and locations based on the EN3 data set (Ingleby and Huddleston,

262 2007), which includes observations from the California Cooperative Fisheries  
263 Investigations (CalCOFI) surveys, as well as Argo and glider data within  
264 our model domain. For biological assimilation, only surface phytoplankton  
265 is used, analogous to what might be obtained under cloud-free conditions  
266 from satellite ocean color data. We note that our data collection for surface  
267 fields is larger than occurs in nature (approximately 13% data coverage on  
268 our model domain in the year 2000), but allows investigation of a best-case,  
269 data-rich scenario. Assimilation of real data is performed in Song et al.  
270 (2016b)

271 Observation errors are added to the sampled data. Errors for physical  
272 variables are drawn from normal random distributions ( $\mathcal{N}(0, 0.1^2)$ ,  $\mathcal{N}(0, 0.01^2)$ ,  
273  $\mathcal{N}(0, 0.02^2)$  and  $\mathcal{N}(0, 0.1^2)$  for in situ temperature, salinity, SSH and SST,  
274 respectively). The observational error levels for in situ temperature, salin-  
275 ity and sea surface height were adopted from Broquet et al. (2009), where  
276 the same data assimilation system was used to fit the data in the same  
277 domain. The observational error level for SST is close to the global stan-  
278 dard deviation of errors (0.13K for AVHRR) (O’Carroll et al., 2012). Errors  
279 in phytoplankton biomass data were drawn randomly from  $\mathcal{N}(0, 0.2^2)$  and  
280 added in log-transformed space. This distribution approximately corresponds  
281 to a 20% multiplicative error which is lower than uncertainty estimates for  
282 global chlorophyll data (Gregg and Casey, 2004, Moore et al., 2009). We also  
283 consider sensitivity experiments in which the observational error for SST is  
284 elevated to 0.4°C and for phytoplankton is increased to 35% and 50%.

285 *3.3. Assimilation setup*

286 Following the method presented in Weaver and Courtier (2001), the back-  
287 ground error covariance is factorized as  $\mathbf{B} = \mathbf{\Sigma}\mathbf{C}\mathbf{\Sigma}^T$ , where  $\mathbf{\Sigma}$  is a diagonal  
288 matrix whose diagonal elements are model standard deviations and  $\mathbf{C}$  is  
289 a correlation matrix. Standard deviations are computed from the 4-year  
290 forward simulation. The background error covariance  $\mathbf{B}_L$  is for  $\ln \mathbf{x}$ , and  
291 therefore biological variables should be log-transformed before computing  
292 the standard deviation. The correlation matrix  $\mathbf{C}$  is obtained through so-  
293 lution of a diffusion equation (Weaver and Courtier, 2001), and we apply  
294 horizontal and vertical length scales of 50 km and 30 m, respectively. It is  
295 reasonable to expect that in general physical and biological variables have  
296 different decorrelation length scales; for example, Lagrangian measurements  
297 in offshore portions of the CCS reveal different decorrelation time-scales for  
298 chlorophyll-a and temperature by Abbott and Letelier (1998). In this study,  
299 we assume that the length scales are identical. Song et al. (2016b) discuss  
300 the requirements of a smaller vertical decorrelation length scale for biological  
301 variables than physical variables in the fully realistic assimilation scenario.

302 The set of experiments proceeds in sequences of 5-day assimilation cycles.  
303 Although Veneziani et al. (2009) has shown that the tangent linear assump-  
304 tion in the physical model is reasonable over a time-scale of 14 days, a shorter  
305 time-scale is required for biological models due to the inherent nonlinearities  
306 of the biological interactions. Song et al. (2016a) find that a time-scale of 5  
307 days is reasonable for the NPZD model and the California Current System  
308 implementation.

309 We examine the coupled assimilative system over the 4-year period 2001–

310 2004, divided into 48, 30-day experiments. Every 30-day experiment consists  
 311 of 6 sequential assimilation cycles, each extending for 5-days. The initial con-  
 312 dition on the first day of each experiment is the 4-year mean state obtained  
 313 for that particular day obtained from the true run. Within each experiment,  
 314 the state estimate at the end of one cycle is used to initialize the background  
 315 estimate for the next 5-day cycle. In our analysis, the first cycle of each ex-  
 316 periment is treated as a spin-up cycle when the linear approximation is the  
 317 least accurate (Song et al., 2016a) and not included in the statistical results.

318 Although we recognize that cross-covariances between model variables  
 319 exist, we calculate univariate correlations only, and we assume that obser-  
 320 vation errors are independent and uncorrelated. These assumptions simplify  
 321 the construction of error covariances  $\mathbf{B}$  and  $\mathbf{R}_i$ , with

$$\mathbf{B} = \begin{bmatrix} \mathbf{B}_G & \mathbf{0} \\ \mathbf{0} & \mathbf{B}_L \end{bmatrix} \quad (5)$$

322 and

$$\mathbf{R}_i = \begin{bmatrix} \mathbf{R}_{G,i} & \mathbf{0} \\ \mathbf{0} & \mathbf{R}_{L,i} \end{bmatrix}. \quad (6)$$

323 For simplicity, we take this highly simplified approach but acknowledge that  
 324 accurate cross-correlations between model variables should yield additional  
 325 improvements in coupled biological and physical state estimation.

326 The reformulation of lognormal 4DVar to the quadratic cost function cre-  
 327 ates an additional linearization approximation. The logarithmic transforma-  
 328 tion of the model state in observation space is linearized using Taylor series  
 329 whose necessary condition is that the difference between the model state and  
 330 observations is small. In order to satisfy this condition, we filter observations



331 such that those more than twice the modeled state are excluded from the as-  
332 simulation procedure (Song et al., 2016a). While this procedure reduces the  
333 total number of assimilated observations it helps ensure the stability of the  
334 data assimilation calculations.

335 We also perform sensitivity experiments for the year 2001 to understand  
336 how increased observational error and the ability of the model to represent  
337 truth impact results. We consider experiments in which either or both SST  
338 and phytoplankton observational error are increased from their original val-  
339 ues. In addition, the coupled data assimilation system is examined after  
340 purposely introducing model error. Physical model error is introduced by  
341 applying surface forcing from the year 2002, and biological model error is  
342 obtained by applying different biological parameter values than the reference  
343 run. Those parameters are listed in Table 1. Physical and biological model  
344 errors are included in the assimilative runs either separately or jointly.

## 345 **4. Performance of the coupled data assimilation system**

### 346 *4.1. The improvement in the Root-mean-squared error*

347 Performance is first evaluated using the root-mean-squared error (RMSE).  
348 RMSEs for surface physical variables – zonal velocity ( $u$ ), meridional veloc-  
349 ity ( $v$ ), SST, and SSH – are significantly improved by coupled assimilation  
350 of physical variables (PDA) (Figure 2). Using the forward model provides  
351 better prediction skill than persistence, but assimilating physical observa-  
352 tions further decreases RMSE. The mean error reduction from the free run  
353 for physical variables is approximately 56% with improvement in both ob-  
354 served variables (SST and SSH) and unobserved variables ( $u$  and  $v$ ). RMSE

355 reduction is larger for observed variables than in unobserved variables, as ex-  
356 pected. Coupled assimilation of biological data (BDA) also decreases RMSE  
357 for physical variables by about 13% from the free run on average, indicat-  
358 ing that surface phytoplankton observations can improve initial conditions  
359 for current, SSH and SST entirely through the dynamics of the coupled sys-  
360 tem. Tracers are influenced by advection and diffusion, and because the  
361 adjoint model includes these coupled dynamics, cost function reduction can  
362 be achieved through alteration of physical as well as biological variables.

363 RMSE reduction is greatest, approximately 57% on average with respect  
364 to the free run, when assimilating both physical and biological data (PBDA).  
365 The modest improvement over PDA indicates that the physical fields are  
366 most constrained by the physical data provided in this experiment. It is  
367 noted that PBDA does not always lead to the improvement over PDA. The  
368 RMSE of PDA for SSH is slightly lower than what PBDA offers, but within  
369 the error bars ( $0.980 \pm 0.017$  cm versus  $0.985 \pm 0.014$  cm). In a case of SST,  
370 PBDA has lower RMSE ( $0.310 \pm 0.006$  °C) than PDA ( $0.314 \pm 0.006$  °C), but  
371 again within error bars. Although the improvement by PBDA over PDA is  
372 statistically marginal for u and v, it is negligible compared with the RMSE  
373 reduction from the free run to PDA.

374 PDA and PBDA reduce the RMSEs for the observed variables (SST and  
375 SSH) more than half almost everywhere (Figure 3(c,d,k,l)). Although the  
376 reduction of RMSE by BDA is less than 20%, BDA reduces the RMSE of  
377 SST and SSH over the entire region (Figure 3(g,h)). Surface currents are  
378 not assimilated, and RMSE reduction in those variables is smaller than for  
379 SST or SSH. There are regions even with the increased RMSE. For instance,

380 PDA shows approximately 20% higher RMSE for  $u$  near the coast between  
381  $40^\circ\text{N}$  and  $45^\circ\text{N}$  (Figure 3(a)). Although the RMSE reduction for  $u$  and  $v$  by  
382 BDA is smaller than for PDA, it occurs over the all model domain (Figure  
383 3(e,f)). PBDA yields the best estimate of surface currents, resulting in the  
384 smallest RMSE. Interestingly, the increased RMSE for  $u$  by PDA in Fig-  
385 ure 3(a) became less obvious when assimilating both physical and biological  
386 observations.

387 Forward simulation of the NPZD model improves the estimation of bi-  
388 ological variables as the RMSEs of the free run are smaller than the result  
389 from one-month persistence. Assimilating data provides a better estimate  
390 for biological variables, further reducing the RMSEs for biological variables  
391 with respect to the free run. In BDA, incorporating surface phytoplank-  
392 ton data improves the estimate of not only phytoplankton but also other  
393 unobserved biological variables. RMSE reduction of biological variables is  
394 approximately 43% on average with respect to that by the free run. PDA  
395 also improves the estimate of biological variables by approximately 7%, and  
396 this improvement results not from the adjustment of biological initial con-  
397 ditions but from the improved representation of the circulation fields. As  
398 described above, biological variables are passively coupled to the currents,  
399 and therefore the sensitivity of the misfits in physical variables to variations  
400 in the biological variables is zero. Improvements in  $T$ ,  $S$  and  $SSH$  cannot be  
401 reduced by changes in biological variables at the initial time. As in the eval-  
402 uation of physical variables, greatest RMSE reduction ( $\sim 49\%$  on average)  
403 for biological variables occurs through coupled assimilation of both physical  
404 and biological data (PBDA), and this reduction is statistically significant.

405 Assimilating physical variables leads to mixed effects on the biological  
406 state estimation at the surface. Figure 5(a-d) show that PDA has a pos-  
407 itive effect near the coastal regions but generally degrades the biological  
408 estimation offshore. Changes in RMSEs by PDA are similar in overall mag-  
409 nitude for all four biological variables. BDA and PBDA result in comparable  
410 RMSE reduction (Figure 5(e-l)). Largest RMSE improvement using these  
411 two methods occurs for phytoplankton, the observed variable. The second  
412 largest reduction in RMSE is seen in detritus. RMSE reduction for P and  
413 D occurs throughout the model domain. The improvements in zooplankton  
414 by BDA and PBDA occur mainly near the coast (Figure 5(f,j)). Least im-  
415 provement is found in the nutrient estimation, and it is visually similar to  
416 the improvement by PDA (Figure 5(c,g,k)).

417 We note that although overall error decreases in all variables, there are  
418 limited regions where the RMSE increases after assimilating surface phyto-  
419 plankton, even in PBDA. Such error increases occur most visibly in Z and  
420 N, but also at one location in D (Figure 5(f,g,j,k)). The NPZD model is a  
421 simple but highly nonlinear system, sometimes stretching the linear approx-  
422 imation used in 4DVar systems. In such cases, the increments can degrade  
423 the posterior estimate. Although this is not limited to our assimilation sys-  
424 tem, it is possible that degrading increments can be amplified due to the  
425 transformation back to the original space using the exponential function.  
426 We note that the limited areas of degradation occur for unobserved vari-  
427 ables only, indicating that the system improves the phytoplankton estimates  
428 through occasionally unrealistic changes to variables for which we have no  
429 information other than background error statistics. Such performance is not

430 surprising in an estimation system, and generally could be improved through  
431 observation of other ecosystem elements.

432 Two additional experiments outlined above better illustrate how modeled  
433 dynamics in coupled data assimilation influence the final state estimate. In  
434 BDAb, biological data is assimilated but only biological variables are ad-  
435 justed. In this case, RMSE in biological fields is reduced by approximately  
436 35% on average (not shown), less than the reduction by BDA with adjust-  
437 ments to all variables. In BDA, the coupled data assimilation system parti-  
438 tions improvements to both physical and biological variables. Adjusting only  
439 biological fields limits the quantitative improvement in biological fields over  
440 the full assimilation cycle relative to what can be achieved through adjust-  
441 ment also of physical fields. RMSE reduction in physical variables is zero in  
442 BDAb because adjusted biological variables do not affect the physical fields.

443 In BDAp physical variables only are adjusted, and misfits in biological  
444 variables can be reduced only through improvement in circulation and mix-  
445 ing. Although the phytoplankton biomass RMSE is reduced in this case,  
446 assimilation generally degrades the estimates of other variables (not shown).  
447 In particular, physical variables are adjusted such that their RMSEs are  
448 larger than in the free run. While the misfit in observed variables can be re-  
449 duced through modification of various fields through coupled dynamics, not  
450 all adjustments result in a better overall estimate of the unobserved variables.

#### 451 *4.2. The improvement in the statistical states*

452 Model performance can be visualized also through a Taylor diagram which  
453 summarizes the model variability relative to the truth, specifically the corre-  
454 lation coefficient and standard deviation (Taylor, 2001). PDA (square termi-

455 nator) and PBDA (circular terminator) show substantial improvements in all  
456 three statistics relative to the prior, especially in observed variables (SST and  
457 SSH) (Figure 6). Unobserved  $u$  and  $v$  are also statistically improved when  
458 physical data are assimilated. For both PDA and PBDA, physical variables  
459 show about the same amount of variability as the true state, and posterior  
460 correlation coefficients are greater than 0.8. Statistically, the improvements  
461 realized by the PDA and PBDA are comparable, implying that physical ob-  
462 servations in this case provide sufficient information for the optimal physical  
463 solution. In BDA (triangular terminator), the statistics of the physical vari-  
464 ables are also improved, with the posterior located closer to the reference  
465 point than the line origins for  $u$ ,  $v$  and SSH, although the improvements  
466 are not as large as for PDA and PBDA. Chlorophyll and ocean currents are  
467 strongly coupled in the advection/diffusion equation while temperature does  
468 not appear in the equations for NPZD model. Hence, chlorophyll observa-  
469 tions have a more substantial impact on  $u, v$  and SSH.

470 In the right panel in Figure 6, results from BDA and PBDA indicate im-  
471 provements in the statistics for all biological variables. As expected, phyto-  
472 plankton, the observed variable, exhibits the greatest improvement, showing  
473 a normalized standard deviation close to 1 and correlation coefficient greater  
474 than 0.9. The statistics of the unobserved variables are also improved, al-  
475 though not as much as for  $P$ . Improvements in the biological error statistics  
476 from BPDA are greater than those for BDA, and more so than improvements  
477 in the physical error statistics by PBDA over PDA; this result indicates again  
478 the significant role that the physical state has on biological fields but not the  
479 reverse. Results from PDA show the least improvement in biological vari-

480 ables.

### 481 *4.3. Sensitivity to observational errors*

482 The performance of data assimilation is dependent on the error levels. We  
483 conduct a sensitivity test with varied observational error levels for SST and  
484 phytoplankton. When observational error for SST is increased from  $0.1^{\circ}\text{C}$   
485 to  $0.4^{\circ}\text{C}$ , the RMSE for SST is elevated by approximately 50% (PDA(0.1,  
486 N/A) versus PDA(0.4, N/A) in Figure 7). Higher SST observational error  
487 also influences the estimation of surface currents and SSH, increasing RM-  
488 SEs more than 20%. The RMSEs for physical variable are less sensitive to  
489 the phytoplankton observational error (BDA(N/A, 20%) versus BDA(N/A,  
490 35%) versus BDA(N/A, 50%) in Figure 7). BDA reduces the RMSEs for  
491 physical variables with respect to the free run even when the phytoplankton  
492 observational error level is 50%. The sensitivity of PBDA to the observa-  
493 tional error can be considered as the mixed response of PDA and BDA to  
494 the changes in observational error for SST and phytoplankton, respectively  
495 (PBDA(0.1, 20%) versus PBDA(0.4, 35%) versus PBDA(0.4, 50%) in Figure  
496 7).

497 Changing observational error for SST does not provide a statistically  
498 significant impact on the RMESs for biological variables (PDA(0.1, N/A)  
499 versus PDA(0.4, N/A) in Figure 8). Higher phytoplankton observational  
500 error degrades the estimation of phytoplankton by elevating the RMSE by  
501 more than 15%, which is statistically significant (BDA(N/A, 20%) versus  
502 BDA(N/A, 35%) versus BDA(N/A, 50%) in Figure 8). Other biological  
503 variables do not have the influence of higher phytoplankton observational  
504 error as much as phytoplankton, showing less than 10% RMSE increase. As

505 for physical variables, the sensitivity of PBDA to the observational error  
506 can also be viewed as the mixed response of PDA and BDA to the changes  
507 in observational error for SST and phytoplankton (PBDA(0.1, 20%) versus  
508 PBDA(0.4, 35%) versus PBDA(0.4, 50%) in Figure 8). It is noted that higher  
509 phytoplankton observational error does not always degrade the estimation for  
510 the biological variables. For instance, BDA(N/A, 50%) shows smaller RMSE  
511 for P than BDA(N/A, 35%). We attribute this to the observation filtering  
512 process discussed in subsection 3.3. More outlying observations are rejected  
513 in BDA(N/A, 50%) than in BDA(N/A, 35%), and eventually the filter helps  
514 the data assimilation run fit observations better, leading to smaller RMSE.

#### 515 *4.4. Sensitivity to the model errors*

516 The performance of coupled physical and biological data assimilation is  
517 also evaluated under the presence of model errors. In this sensitivity test,  
518 the control run is PBDA with 0.4°C and 35% observational errors for SST  
519 and phytoplankton, respectively. We refer to EF as error free, EP as error in  
520 physics (in which the wrong year’s surface forcing has been introduced), EB  
521 as error in biology (with different NPZD model parameters than the reference  
522 run), and EPB as error in physics and biology. PBDA results in higher  
523 RMSEs in the observed variables, SST and SSH, when the surface forcing of  
524 the year 2002 is used in the assimilative run for the year 2001, but RMSEs  
525 are still considerably lower than that of the free run in which no model  
526 error is included (PBDA, EF versus PBDA, EP in Figure 9). However, the  
527 introduced physical model error degrades the estimation of surface currents:  
528 RMSEs of PBDA, EP for u and v are greater than that of the free run with  
529 no model error. The impact of biological model error to physical variables



530 is not statistically significant. Using incorrect biological parameter values in  
531 PBDA does not change RMSEs for all physical variables (PBDA, EP versus  
532 PBDA, EPB in Figure 9).

533 Estimating the biological state can be influenced by both physical and  
534 biological model errors. Using the wrong forcing degrades the estimate of  
535 biological variables, leading higher RMSEs (PBDA, EF versus PBDA, EP in  
536 Figure 10). The degradation is particularly strong for N and D, resulting in  
537 higher RMSEs than those from the free run. The introduced biological model  
538 error also has statistically significant impact on the estimation of biological  
539 variables. Using wrong parameter values and surface forcing increased the  
540 RMSE for phytoplankton, although it is still slightly smaller than that from  
541 the free run with no model error (PBDA, EF versus PBDA, EP versus PBDA,  
542 EPB in Figure 10). However, our model errors result in higher RMSEs for  
543 unobserved biological variables. Including both physical and biological errors  
544 makes the RMSEs for Z, N and D greater with respect to the free run. The  
545 most substantial impact is observed in Z and is perhaps results from the fact  
546 that two of four modified parameters are associated with the zooplankton  
547 dynamics.

548 The incremental form of 4DVar used here does not allow easily for the  
549 accurate computation of posterior error estimates. The posterior error co-  
550 variance is equivalent to the inverse of Hessian matrix (Moore et al., 2012).  
551 Here, the inverse of Hessian matrix is estimated as  $\mathbf{V}\mathbf{T}^{-1}\mathbf{V}^T$ , where  $\mathbf{V}$  is the  
552 orthogonal matrix with Lanczos vectors, and  $\mathbf{T}$  is the symmetric tridiagonal  
553 matrix that contains coefficients in Lanczos recurrence relation (Song et al.,  
554 2016a). The Hessian matrix is approximated by a tridiagonal factorization

555 using the Lanczos vectors, and represents a reduced approximation. The  
556 leading eigenvectors of the Hessian matrix can be estimated from the Lanc-  
557 zos vectors but these would represent the smallest eigenvectors of the analysis  
558 error covariance. Therefore, this calculation provides a poor representation of  
559 the analysis error covariance matrix. One can consider other effective meth-  
560 ods (e.g. Daescu and Navon, 2007, a reduced second order adjoint model)  
561 for the estimation of Hessian, but the inverse of the Hessian matrix is still  
562 approximated with the least important orthogonal vectors. We note that  
563 that ROMS does provide options to estimate the posterior error covariance  
564 in dual form (Moore et al., 2012), but that is not the form used in this study.

## 565 **5. Summary and Discussion**

566 We have developed and investigated combined physical and biological 4-  
567 dimensional variational data assimilation in an ocean model. Biological data  
568 assimilation benefits from a unique approach because of the non-Gaussian  
569 statistics of biological variables and their errors. We have assumed lognor-  
570 mal statistics for these variables and applied the quadratic formulation of the  
571 incremental approximation developed by Song et al. (2016a). Assimilation  
572 of variables having different error statistics is required for combined physi-  
573 cal and biological assimilation and proceeds here following the approach of  
574 Fletcher (2010) and Fletcher and Jones (2014).

575 In model twin experiments using ROMS and a 4-component, NPZD  
576 ecosystem model configured for the realistic California Current System, we  
577 investigated how coupled biological and physical data assimilation improves  
578 overall estimates of the combined physical and biological ocean state. Ob-

579 observations were drawn from a forward model run which represented the true  
580 ocean state with errors drawn from distributions with known statistics added.  
581 Three observation sampling strategies were chosen: (1) biological observa-  
582 tions, (2) physical observations, and (3) both biological and physical ob-  
583 servations. Then we altered the model state through adjustment of initial  
584 conditions in physical and/or biological variables using those observations.  
585 Statistics of RMSE, correlation coefficients, state variability were analyzed  
586 from a total of 48 sequences of six 5-day assimilation cycles.

587 We found that assimilation of physical data (PDA) improves model-data  
588 misfit of physical variables, and assimilation of biological data (BDA) re-  
589 duces model error for biological variables. Such results should be expected.  
590 In addition, PDA resulted in biological error reduction, and BDA yielded  
591 improvements in the physical variable misfit. Even though PDA has no ef-  
592 fect on biological initial conditions, it does influence biological variables over  
593 the entire assimilation cycle through improvements to the physical fields  
594 which then feed back on biological variables through tangent linear (and  
595 forward) model advection and diffusion. In contrast, a change in biological  
596 initial conditions, for example resulting from BDA, has no influence on phys-  
597 ical variables through the forward model directly. However, BDA influences  
598 physical variable initial conditions through coupled dynamics included in the  
599 adjoint model. The coupled data assimilation system partitions changes to  
600 all control vector elements that can reduce error in the biological data misfit,  
601 and this partitioning extends to both biological and physical variables. This  
602 conceptual division of influence is drawn schematically in Figure 11.

603 Overall, the greatest performance in both physical and biological fields

604 as quantified by various measures resulted from the combined assimilation of  
605 physical and biological data (PBDA), further supporting the interpretation  
606 of the PDA and BDA results. While physical observations provide the most  
607 effective constraint for physical variables, and biological observations most  
608 constrain biological variables, additional improvement in physical variables  
609 derived from biological information through model adjoint dynamics and bi-  
610 ological errors can be reduced through physical observations via the tangent  
611 linear (and forward) model. Higher observational errors in SST and phyto-  
612 plankton increase the RMSEs of PBDA, but the increase is smaller than the  
613 combined increments in RMSE of PDA with higher SST observational er-  
614 ror and of BDA with higher phytoplankton observational error. Introducing  
615 model errors (through the application of incorrect surface forcing or altered  
616 biological parameters) also degrades the performance of PBDA, but its im-  
617 pact on PBDA is not particularly different from that for PDA and BDA on  
618 a monthly time scale.

619 More generally, variables in a coupled 4DVar system can be influenced  
620 in two ways, dynamically through the adjoint and tangent linear models  
621 and statistically through covariances of the background error covariance ma-  
622 trix. In this study, univariate spatial correlations in fields were assumed  
623 through the integration of a diffusion equation; however, no multivariate  
624 correlations were represented, and therefore all improvements resulting from  
625 observations of coupled variables resulted exclusively from adjoint and tan-  
626 gent linear model dynamics. These coupling dynamics are not included in  
627 alternate data assimilation approaches based on statistical estimation alone,  
628 and statistical correlations between physical and biological variables provide

629 the only way to transfer this critical information. These correlations are not  
630 well-known in nature, though some groups have reported such information  
631 (Behrenfeld et al., 2006), and estimates that are consistent with the ocean  
632 circulation and in principle ecosystem models can be calculated from forward  
633 model calculations (Shulman et al., 2013). Ensemble-based data assimilation  
634 approaches use the ensemble to estimate time-dependent correlations (Simon  
635 et al., 2015), usually with an inflation factor and localization to compensate  
636 the effects of having a small ensemble. We would expect that the present  
637 4DVar assimilation approach would further benefit from better background  
638 error covariance estimates, a subject for future study.

639 The variational approach to coupled dynamics with mixed statistics pre-  
640 sented here is conceptually straightforward to implement within any existing  
641 coupled system equipped with tangent and adjoint models and assuming no  
642 multivariate correlations between physical and biological variables. The com-  
643 putational cost of the combined physical and biological system is comparable  
644 to the cost of either the physical or biological system in isolation. Results  
645 from this study suggest that coupled assimilation using 4DVar is practical  
646 and realizable. However, the twin experiment framework used here repre-  
647 sents an idealized setting in which the data is unencumbered by cloud cover  
648 and the model surface forcing and boundary conditions are error free.

649 Our conclusions are drawn based on the ensemble of 30-day assimilation  
650 (6 cycles). It is possible that the model may drift from the truth if biases  
651 are introduced by assimilation and accumulate over time scales longer than  
652 one month. In our ideal twin experiment, error is introduced in the initial  
653 conditions only and estimating accurate initial conditions always improves

654 the model bias. Indeed, the model bias does not increase for more than 6  
655 months when we tested the data assimilation system using 48 cycles (not  
656 shown). However, in realistic scenarios in which errors in surface forcing,  
657 boundary forcing, and model construction exist, it is possible that model  
658 bias develops more rapidly. Such an issue does not appear in Song et al.  
659 (2016b), but further studies of model bias in realistic scenarios over long  
660 periods of time is warranted.

661 When assimilating real observations, the presented assimilation system  
662 may encounter obstacles. For instance, the model dynamics inevitably mis-  
663 represent or entirely miss important processes in nature. Under these circum-  
664 stances, adjustments to the initial conditions determined by model dynamics  
665 are of limited value in matching observations. Large differences between the  
666 observed and prior state can also create an issue because they violate the  
667 linearization of log-transformation function and may prevent solution con-  
668 vergence. In order to prevent this outcome, the coupled assimilation system  
669 requires a filtering process that excludes observations far from the prior. The  
670 filtering process may reduce the number of observations, but it stabilizes the  
671 assimilation system and may lead to a better posterior solution as shown  
672 in subsection 4.3. We note that the filtering procedure is reevaluated during  
673 each outer loop of the assimilation system, and observations that are rejected  
674 initially may be included in the final outer loop.

675 In a companion paper Song et al. (2016b), we investigate the assimila-  
676 tion system's performance in a more realistic system in which real remotely  
677 sensed and in situ physical and ecosystem data are assimilated. In that real-  
678 istic setting, improvements to RMSE for physical variables is not improved

679 by assimilating real chlorophyll observations. Several factors may account  
680 for this result that stands in contrast to that offered by the twin experi-  
681 ments here. Data availability is reduced relative to this study as frequent  
682 cloud cover prevents collection of SST and chlorophyll data over much of our  
683 domain. In addition, physical and biological model error are likely greater  
684 than that considered in the present study. Even though the real assimila-  
685 tion experiment is carried out on a higher resolution grid that better resolves  
686 the CCS mesoscale circulation, the physical model is still imperfect relative  
687 to nature. The NPZD ecosystem model used is advantageous for its rela-  
688 tive simplicity, and has been applied to multiple realistic studies of ocean  
689 biogeochemistry, including in the CCS (Powell et al., 2006). However, with  
690 only one phytoplankton functional group, it is less than ideal in representing  
691 the multiple phytoplankton communities in different geographical regions of  
692 the CCS. Assimilation improvements may result from application of a more  
693 complex biogeochemical model. At this time, we do not know which of  
694 these elements is responsible for the differences between the two studies, but  
695 the present study shows that under excellent conditions in which a model  
696 is nearly able to represent truth and observations are abundant, the lowest  
697 RMSE for physical and biological variables results from assimilation of both  
698 biological and physical variables into the coupled system.

## 699 **6. Acknowledgement**

700 We are grateful for grants from the Gordon and Betty Moore Foundation  
701 and from the National Oceanographic and Atmospheric Administration office  
702 of Oceanic and Atmospheric Research (award number NA10OAR4320156)

703 that supported this research. The authors would like to thank four anony-  
704 mous reviewers for valuable comments and suggestions, which significantly  
705 improved the manuscript.

706 Abbott, M. R., Letelier, R. M., 1998. Decorrelation scales of chlorophyll  
707 as observed from bio-optical drifters in the california current. *Deep Sea*  
708 *Research Part II: Topical Studies in Oceanography* 45 (89), 1639 – 1667.

709 Anderson, L. A., Robinson, A. R., Lozano, C. J., 2000. Physical and biologi-  
710 cal modeling in the Gulf Stream region: I. Data assimilation methodology.  
711 *Deep Sea Res. Pt I* 47, 1787–1827.

712 Balmaseda, M. A., Vidard, A., Anderson, D. L. T., 2008. The ECMWF  
713 Ocean Analysis System: ORA-S3. *Mon. Wea. Rev.* 136, 3018–3034.

714 Behrenfeld, M. J., O'Malley, R. T., Siegel, D. A., McClain, C. R., Sarmiento,  
715 J. L., Feldman, G. C., Milligan, A. J., Falkowski, P. G., Letelier, R. M.,  
716 Boss, E. S., 2006. Climate-driven trends in contemporary ocean produc-  
717 tivity. *Nature* 444, 752–755.

718 Behringer, D. W., Ji, M., Leetmaa, A., 1998. An improved coupled model  
719 for ENSO prediction and implications for ocean initialization. Part I: The  
720 ocean data assimilation system. *Mon. Wea. Rev.* 126, 1013–1021.

721 Bell, M. J., Forbes, R. M., Hines, A., 2000. Assessment of the FOAM global  
722 data assimilation system for real-time operational ocean forecasting. *J.*  
723 *Mar. Sys.* 25, 1–22.

724 Berline, L., Brankart, J. M., Brasseur, P., Ourmières, Y., Verron, J., 2007.  
725 Improving the physics of a coupled physical-biogeochemical model of the



- 726 North Atlantic through data assimilation: Impact on the ecosystem. *J.*  
727 *Marine Syst.* 64, 153–172.
- 728 Bertino, L., Evensen, G., Wackernagel, H., 2003. Sequential data assimilation  
729 techniques in oceanography. *Int. Statist. Rev.* 71, 223–241.
- 730 Brink, K. H., Cowles, T. J., 1991. The coastal transition zone program.  
731 *Journal of Geophysical Research: Oceans* 96 (C8), 14637–14647.
- 732 Broquet, G., Edwards, C. A., Moore, A. M., Powell, B. S., Veneziani, M.,  
733 Doyle, J. D., 2009. Application of 4D-Variational data assimilation to the  
734 California Current System. *Dynam. Atmos. Oceans* 48, 69–92.
- 735 Broquet, G., Moore, A. M., Arango, H. G., Edwards, C. A., 2011. Correc-  
736 tions to ocean surface forcing in the California Current System using 4D  
737 variational data assimilation. *Ocean Modell.* 36, 116–132.
- 738 Campbell, J. W., 1995. The lognormal distribution as a model for bio-optical  
739 variability in the sea. *J. Geophys. Res.* 100 (C7), 13237–13254.
- 740 Carton, J., Giese, B., 2008. A reanalysis of ocean climate using Simple Ocean  
741 Data Assimilation (SODA). *Mon. Wea. Rev.* 136, 2999–3017.
- 742 Chai, F., Liu, G., Xue, H., Shi, L., Chao, Y., Tseng, C.-M., Chou, W.-C.,  
743 Liu, K.-K., 2009. Seasonal and interannual variability of carbon cycle in  
744 South China Sea: A three-dimensional physical-biogeochemical modeling  
745 study. *Journal of Oceanography* 65 (5), 703–720.
- 746 Chassignet, E. P., Hurlburt, H. E., Smedstad, O. M., Halliwell, G. R., Hogan,  
747 P. J., Wallcraft, A. J., Baraille, R., Bleck, R., 2007. The HYCOM (HYbrid

- 748 Coordinate Ocean Model) data assimilative system. *J. Marine Syst.* 65,  
749 60–83.
- 750 Chenillat, F., Rivière, P., Capet, X., Franks, P. J. S., Blanke, B., 2013.  
751 California coastal upwelling onset variability: Cross-shore and bottom-up  
752 propagation in the planktonic ecosystem. *PLoS ONE* 8 (5).
- 753 Ciavatta, S., Torres, R., Saux-Picart, S., Allen, J. I., 2011. Can ocean color  
754 assimilation improve biogeochemical hindcasts in shelf seas? *J. of Geophys.*  
755 *Res.: Oceans* 116 (C12), n/a–n/a.
- 756 Cox, P. M., Betts, R. A., Spall, S. A., Totterdell, I. J., 2000. Acceleration of  
757 global warming due to carbon-cycle feedbacks in a coupled climate model.  
758 *Nature* 408, 184–187.
- 759 Cummings, J., Bertino, L., Brasseur, P., Fukumori, I., Kamachi, M., Martin,  
760 M., Mogensen, K., Oke, P., Testut, C., Verron, J., Weaver, A., 2009. Ocean  
761 data assimilation systems for GODAE. *Oceanography* 22 (3), 96–109.
- 762 Daescu, D., Navon, I., 2007. Efficiency of a POD-based reduced second order  
763 adjoint model in 4D-Var data assimilation. *Int. J. Numer. Methods Fluids*  
764 53, 985–1004.
- 765 Doron, M., Brasseur, P., Brankart, J.-M., Losa, S., Mellet, A., 2013.  
766 Stochastic estimation of biogeochemical parameters from Globcolour ocean  
767 colour satellite data in a North Atlantic 3D ocean coupled physical-  
768 biogeochemical model. *J. Marine Syst.* 117–118, 81–95.

- 769 Doyle, J. D., Jiang, Q., Chao, Y., Farrara, J., 2009. High-resolution real-  
770 time modeling of the marine atmospheric boundary layer in support of the  
771 AOSN-II field campaign. *Deep-Sea Res. Pt. II* 56, 87–99.
- 772 Edwards, C. A., Moore, A. M., Hoteit, I., Cornuelle, B. D., 2015. Regional  
773 ocean data assimilation. *Annu. Rev. Mar. Sci.* 7, 6.1–6.22.
- 774 Fennel, K., 2010. The role of continental shelves in nitrogen and carbon  
775 cycling: Northwestern North Atlantic case study. *Ocean Sci.* 6, 539–548.
- 776 Fiechter, J., Broquet, G., Moore, A. M., Arango, H. G., 2011. A data assim-  
777 ilative, coupled physical-biological model for the coastal Gulf of Alaska.  
778 *Dynam. Atmos. Oceans* 51, 75–98.
- 779 Fletcher, S. J., 2010. Mixed Gaussian-lognormal four-dimensional data as-  
780 similation. *Tellus A* 62, 266–287.
- 781 Fletcher, S. J., Jones, A. S., 2014. Multiplicative and additive incremental  
782 variational data assimilation for mixed lognormal-gaussian errors. *Mon.*  
783 *Wea. Rev.* 142, 2521–2544.
- 784 Fletcher, S. J., Zupanski, M., 2006. A hybrid multivariate normal and log-  
785 normal distribution for data assimilation. *Atmosph. Sci. Lett.* 7, 43–46.
- 786 Follows, M. J., Dutkiewicz, S., Grant, S., Chisholm, S. W., 2007. Emer-  
787 gent biogeography of microbial communities in a model ocean. *Science*  
788 315 (5820), 1843–1846.
- 789 Ford, D. A., Edwards, K. P., Lea, D., Barciela, R. M., Martin, M. J., De-

- 790 maria, J., 2012. Assimilating globcolour ocean colour data into a pre-  
791 operational physical-biogeochemical model. *Ocean Science* 8 (5), 751–771.
- 792 Franks, P. J. S., Chen, C., 2001. A 3-D prognostic numerical model study of  
793 the Georges bank ecosystem. Part ii: biologicalphysical model. *Deep-Sea*  
794 *Res. II* 48, 457–482.
- 795 Friedrichs, M. A. M., 2001. Assimilation of JGOFS EqPac and SeaWiFS data  
796 into a marine ecosystem model of the central equatorial Pacific ocean. *Deep*  
797 *Sea Res. Pt II* 49, 289–319.
- 798 Frouin, R., Iacobellis, S. F., 2002. Influence of phytoplankton on the global  
799 radiation budget. *J. Geophys. Res.: Atmospheres* 107 (D19), ACL 5–1–  
800 ACL 5–10, 4377.  
801 URL <http://dx.doi.org/10.1029/2001JD000562>
- 802 Garcia-Gorriz, E., Hoepffner, N., Ouberdous, M., 2003. Assimilation of Sea-  
803 WiFS data in a coupled physical-biological model of the Adriatic Sea. *J.*  
804 *Marine Syst.* 40-41, 233–252.
- 805 Goebel, N. L., Edwards, C. A., Zehr, J. P., Follows, M. J., 2010. An emergent  
806 community ecosystem model applied to the California Current System. *J.*  
807 *Marine Syst.* 83.
- 808 Gregg, W. W., 2008. Assimilation of SeaWiFS ocean chlorophyll data into a  
809 three-dimensional global ocean model. *J. Marine Syst.* 69, 205 – 225.
- 810 Gregg, W. W., Casey, N. W., 2004. Global and regional evaluation of the  
811 SeaWiFS chlorophyll data set. *Remote Sens. Environ.* 93, 463 – 479.

- 812 Hoteit, I., Triantafyllou, G., Petihakis, G., 2005. Efficient data assimilation  
813 into a complex, 3-d physical-biogeochemical model using partially-local  
814 Kalman filters. *Ann. Geophys.* 23, 3171–3185.
- 815 Hu, J., Fennel, K., Mattern, J. P., Wilkin, J., 2012. Data assimilation with  
816 a local ensemble Kalman filter applied to a three-dimensional biological  
817 model of the Middle Atlantic Bight. *J. Marine Syst.* 94, 145 – 156.
- 818 Ingleby, B., Huddleston, M., 2007. Quality control of ocean temperature and  
819 salinity profiles - Historical and real time data. *J. Mar. Syst.* 65, 158–175.
- 820 Köhl, A., Stammer, D., Cornuelle, B. D., 2007. Interannual to decadal  
821 changes in the ECCO global synthesis. *J. Phys. Oceanogr.* 37, 313–337.
- 822 Kurapov, A. L., Foley, D., Strub, P. T., Egbert, G. D., Allen, J. S., 2011.  
823 Variational assimilation of satellite observations in a coastal ocean model  
824 off Oregon. *J. Geophys. Res.: Oceans* 116 (C5).
- 825 Lewis, M., Carr, M., Feldman, G., Esaias, W., McClain, C., 1990. Influence  
826 of penetrating solar radiation on the heat budget of the equatorial Pacific  
827 Ocean. *Nature* 347 (6293), 543–545.
- 828 Losa, S., Kivman, G., Ryabchenko, V., 2004. Weak constraint parameter  
829 estimation for a simple ocean ecosystem model: What can we learn about  
830 the model and data? *J. Marine Syst.* 45, 1–20.
- 831 Matear, R. J., 1995. Parameter optimization and analysis of ecosystem mod-  
832 els using simulated annealing: A case study at Station P. *J. Mar. Res.*  
833 53 (4), 571–607.

- 834 Mattern, J., Fennel, K., Dowd, M., 2012. Estimating time-dependent param-  
835 eters for a biological ocean model using an emulator approach. *J. Marine*  
836 *Syst.* 96–97, 32–47.
- 837 Matthews, D., Powell, B. S., Janekovi, I., 2012. Analysis of four-dimensional  
838 variational state estimation of the hawaiian waters. *J. of Geophys. Res.:*  
839 *Oceans* 117 (C3).
- 840 Miller, A. J., Di Lorenzo, E., Neilson, D. J., Cornuelle, B. D., Moisan, J. R.,  
841 2000. Modeling CalCOFI observations during El Niño: Fitting physics and  
842 biology. *Calif. Coop. Ocean. Fish. Invest. Rep.* 41, 87–97.
- 843 Moore, A., Arango, H., Broquet, G., Edwards, C. A., Veneziani, M., Pow-  
844 ell, B., Foley, D., Doyle, J., Costa, D., Robinson, P., 2011a. The Regional  
845 Ocean Modeling System (ROMS) 4-dimensional variational data assimila-  
846 tion systems, Part III: Observation impact and observation sensitivity in  
847 the California Current System. *Prog. Oceanogr.* 91, 74–94.
- 848 Moore, A. M., Arango, H. G., Broquet, G., 2012. Estimates of analysis and  
849 forecast error variances derived from the adjoint of 4D-Var. *Mon. Wea.*  
850 *Rev.* 140, 3183–3203.
- 851 Moore, A. M., Arango, H. G., Broquet, G., Edwards, C. A., Veneziani, M.,  
852 Powell, B. S., Foley, D., Doyle, J., Costa, D., Robinson, P., 2011b. The  
853 Regional Ocean Modeling System (ROMS) 4-dimensional variational data  
854 assimilation systems, Part II: Performance and application to the Califor-  
855 nia Current System. *Prog. Oceanogr.* 91, 50–73.

- 856 Moore, A. M., Arango, H. G., Broquet, G., Powell, B. S., Zavala-Garay,  
857 J., Weaver, A. T., 2011c. The Regional Ocean Modeling System (ROMS)  
858 4-dimensional variational data assimilation systems, Part I: Formulation  
859 and Overview. *Prog. Oceanogr.* 91, 34–49.
- 860 Moore, T. S., Campbell, J. W., Dowell, M. D., 2009. A class-based approach  
861 to characterizing and mapping the uncertainty of the MODIS ocean chloro-  
862 phyll product. *Remote Sens. Environ.* 113 (11), 2424–2430.
- 863 Morel, A., 1988. Optical modeling of the upper ocean in relation to its  
864 biogenous matter content (case i waters). *J. Geophys. Res.: Oceans*  
865 93 (C9), 10749–10768.  
866 URL <http://dx.doi.org/10.1029/JC093iC09p10749>
- 867 Murtugudde, R., Beauchamp, J., McClain, C. R., Lewis, M., Busalacchi,  
868 A. J., 2002. Effects of penetrative radiation on the upper tropical ocean  
869 circulation. *J. Climate* 15 (5), 470–486.
- 870 Natvik, L. J., Evensen, G., 2003. Assimilation of ocean colour data into  
871 a biochemical model of the North Atlantic: Part 1. Data assimilation  
872 experiments. *J. Marine Syst.* 40-41, 127–153.
- 873 O’Carroll, A. G., August, T., Borgne, P. L., Marsouin, A., 2012. The accu-  
874 racy of SST retrievals from Metop-A IASI and AVHRR using the EUMET-  
875 SAT OSI-SAF matchup dataset. *Remote Sens. Environ.* 126, 184–194.
- 876 Oke, P., Brassington, G., Griffin, D., Schiller, A., 2008. The Bluelink ocean  
877 data assimilation system (BODAS). *Ocean Modell.* 21, 46–70.

- 878 Oschlies, A., Garçon, V., 1998. Eddy enhancement of primary production in  
879 a model of the North Atlantic Ocean. *Nature* 394, 266–269.
- 880 Ourmières, Y., Brasseur, P., Lévy, M., Brankart, J.-M., Verron, J., 2009. On  
881 the key role of nutrient data to constrain a coupled physicalbiogeochemical  
882 assimilative model of the North Atlantic Ocean. *J. Marine Syst.* 75, 100–  
883 115.
- 884 Park, J.-Y., Kug, J.-S., Bader, J., Rolph, R., Kwon, M., 2015. Amplified arc-  
885 tic warming by phytoplankton under greenhouse warming. *PNAS* 112 (19),  
886 5921–5926.
- 887 Powell, T., Lewis, C., Curchitser, E., Haidvogel, D., Hermann, A., Dob-  
888 bins, E., 2006. Results from a three-dimensional, nested, biologicalphysi-  
889 cal model of the california current system and comparisons with statistics  
890 from satellite imagery. *J. Geophys. Res.* 111, C07018.
- 891 Raghukumar, K., Edwards, C. A., Goebel, N. L., Broquet, G., Veneziani, M.,  
892 Moore, A. M., Zehr, J. P., 2015. Impact of assimilating physical oceanog-  
893 raphic data on modeled ecosystem dynamics in the California Current  
894 System. *Prog. Oceanogr.* 138 (0), 546–558.
- 895 Rousseaux, C. S., Gregg, W. W., 2012. Climate variability and phyto-  
896 plankton composition in the Pacific Ocean. *J. of Geophys. Res.: Oceans*  
897 117 (C10), n/a–n/a.
- 898 Roy, S., Broomhead, D., Platt, T., Sathyendranath, S., Ciavatta, S., 2012.  
899 Sequential variations of phytoplankton growth and mortality in an NPZ  
900 model: A remote-sensing-based assessment. *J. Marine Syst.* 92, 16–29.



- 901 Sakov, P., Counillon, F., Bertino, L., Lister, K., Oke, P., Korabely, A., 2012.  
902 TOPAZ4: an ocean-sea ice data assimilation system for the North Atlantic  
903 and Arctic. *Ocean Science* 8 (4), 633–656.
- 904 Shulman, I., Frolov, S., Anderson, S., Penta, B., Gould, R., Sakalaukus,  
905 P., Ladner, S., 2013. Impact of bio-optical data assimilation on short-  
906 term coupled physical, bio-optical model predictions. *J. of Geophys. Res.:  
907 Oceans* 118 (4), 2215–2230.
- 908 Shulman, I., Rowley, C., Anderson, S., DeRada, S., Kindle, J., Martin, P.,  
909 Doyle, J., Cummings, J., Ramp, S., Chavez, F., Fratantoni, D., Davis,  
910 R., 2009. Impact of glider data assimilation on the monterey bay model.  
911 *Deep-Sea Res. II* 56 (35), 188 – 198.
- 912 Simon, E., Bertino, L., 2009. Application of the Gaussian anamorphosis to  
913 assimilation in a 3-D coupled physical-ecosystem model of the North At-  
914 lantic with the EnKF: a twin experiment. *Ocean Sci.* 5, 495–510.
- 915 Simon, E., Samuelsen, A., Bertino, L., Mouysset, S., 2015. Experiences in  
916 multiyear combined state-parameter estimation with an ecosystem model  
917 of the North Atlantic and Arctic Oceans using the Ensemble Kalman Filter.  
918 *J. Marine Syst.* 152, 1–17.
- 919 Song, H., Edwards, C. A., Moore, A. M., Fiechter, J., 2012. Incremental  
920 four-dimensional variational data assimilation of positive-definite oceanic  
921 variables using a logarithm transformation. *Ocean Modell.* 54–55, 1–17.
- 922 Song, H., Edwards, C. A., Moore, A. M., Fiechter, J., 2016a. Data assimi-  
923 lation in a coupled physical-biogeochemical model of the California Cur-

924     rent System using an incremental lognormal 4-dimensional variational ap-  
925     proach: Part 1, Model formulation and biological data assimilation twin  
926     experiments. *Ocean Modell.* In Press.

927     Song, H., Edwards, C. A., Moore, A. M., Fiechter, J., 2016b. Data assimi-  
928     lation in a coupled physical-biogeochemical model of the California Cur-  
929     rent System using an incremental lognormal 4-dimensional variational ap-  
930     proach: Part 3, Assimilation in a realistic context using satellite and in  
931     situ observations. *Ocean Modell.* In Press.

932     Spitz, Y., Moisan, J., Abbott, M., Richman, J., 1998. Data assimilation and a  
933     pelagic ecosystem model: parameterization using time series observations.  
934     *J. Mar. Systems* 16 (12), 51 – 68.

935     Stock, C. A., Alexander, M. A., Bond, N. A., Brander, K. M., Cheung,  
936     W. W., Curchitser, E. N., Delworth, T. L., Dunne, J. P., Griffies, S. M.,  
937     Haltuch, M. A., Hare, J. A., Hollowed, A. B., Lehodey, P., Levin, S. A.,  
938     Link, J. S., Rose, K. A., Rykaczewski, R. R., Sarmiento, J. L., Stouffer,  
939     R. J., Schwing, F. B., Vecchi, G. A., Werner, F. E., 2011. On the use  
940     of IPCC-class models to assess the impact of climate on Living Marine  
941     Resources. *Prog. Oceanogr.* 88, 1–27.

942     Taylor, K. E., 2001. Summarizing multiple aspects of model performance in  
943     a single diagram. *J. Geophys. Res.* 106, 7183–7192.

944     Triantafyllou, G., Korres, G., Hoteit, I., Petihakis, G., Banks, A., 2007.  
945     Assimilation of ocean colour data into a biogeochemical flux model of the  
946     eastern Mediterranean Sea. *Ocean Sci.* 3, 397–410.

Table 1: A list of parameters for the NPZD model that is chosen differently in the assimilation runs to include the biological model error. The columns from the left to the right represent names, units, values in the true run and values for biological model error, respectively. The parameter values for the biological model error are from the study for Gulf of Alaska in Fiechter et al. (2011).

Parameter name	unit	Value, True	Value, EB
Uptake rate for nitrate	$day^{-1}$	1.0	0.8
Zooplankton grazing rate	$day^{-1}$	0.65	0.4
Ivlev constant	Dimensionless	0.4	0.84
Detritus remineralization rate	Dimensionless	0.1	0.2

947 Veneziani, M., Edwards, C. A., Moore, A. M., 2009. A central california  
 948 coastal ocean modeling study: 2. adjoint sensitivities to local and remote  
 949 forcing mechanisms. J. Geophys. Res. 114, C04020.

950 Weaver, A., Courtier, P., 2001. Correlation modelling on the sphere using a  
 951 generalized diffusion equation. Quart. J. Roy. Meteorol. Soc. 127, 1815–  
 952 1846.

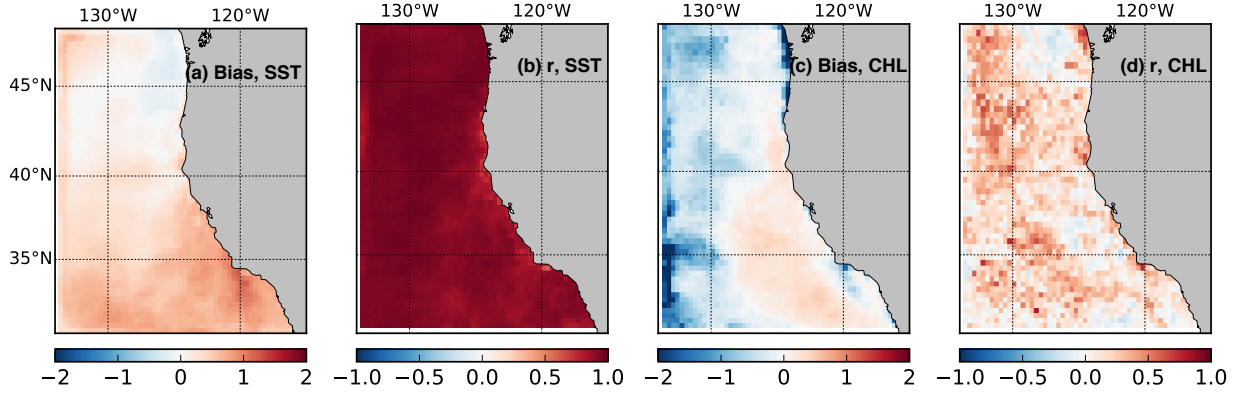


Figure 1: Normalized bias of the simulated (a) SST and (c) surface chlorophyll that are assumed as the truth in a twin experiment, and their correlation coefficients (b, d). Monthly mean values from the model and the satellite observation AVHRR (SeaWiFS), as the reference states for SST (surface chlorophyll), are considered in the calculation.

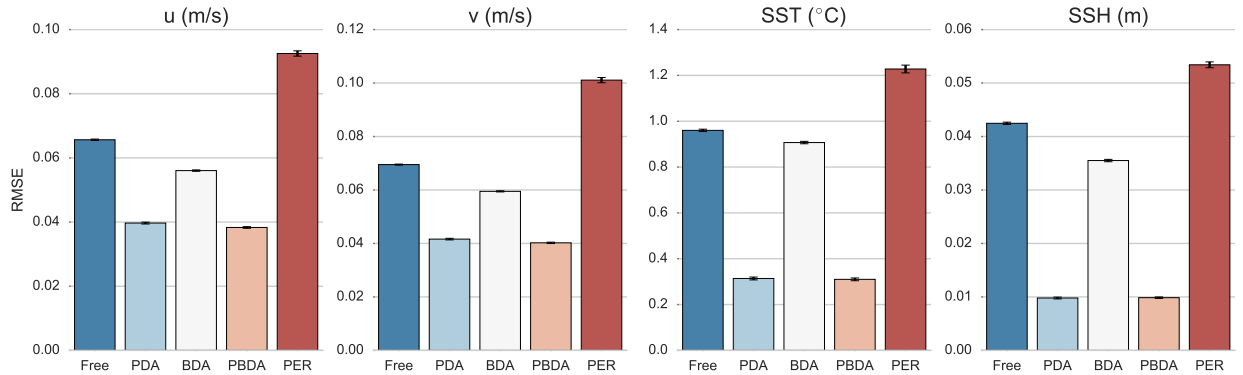


Figure 2: Root-mean-squared error (RMSE) of  $u$ ,  $v$ , SST and SSH at the surface in the four different simulation: Free run (blue), analysis by the PDA (light blue), analysis by the BDA (white), analysis by the PBDA (light red). Red bar represents the one-month persistence RMSE. Error bars represent standard error from 1200 days (25 days × 12 months × 4 years).

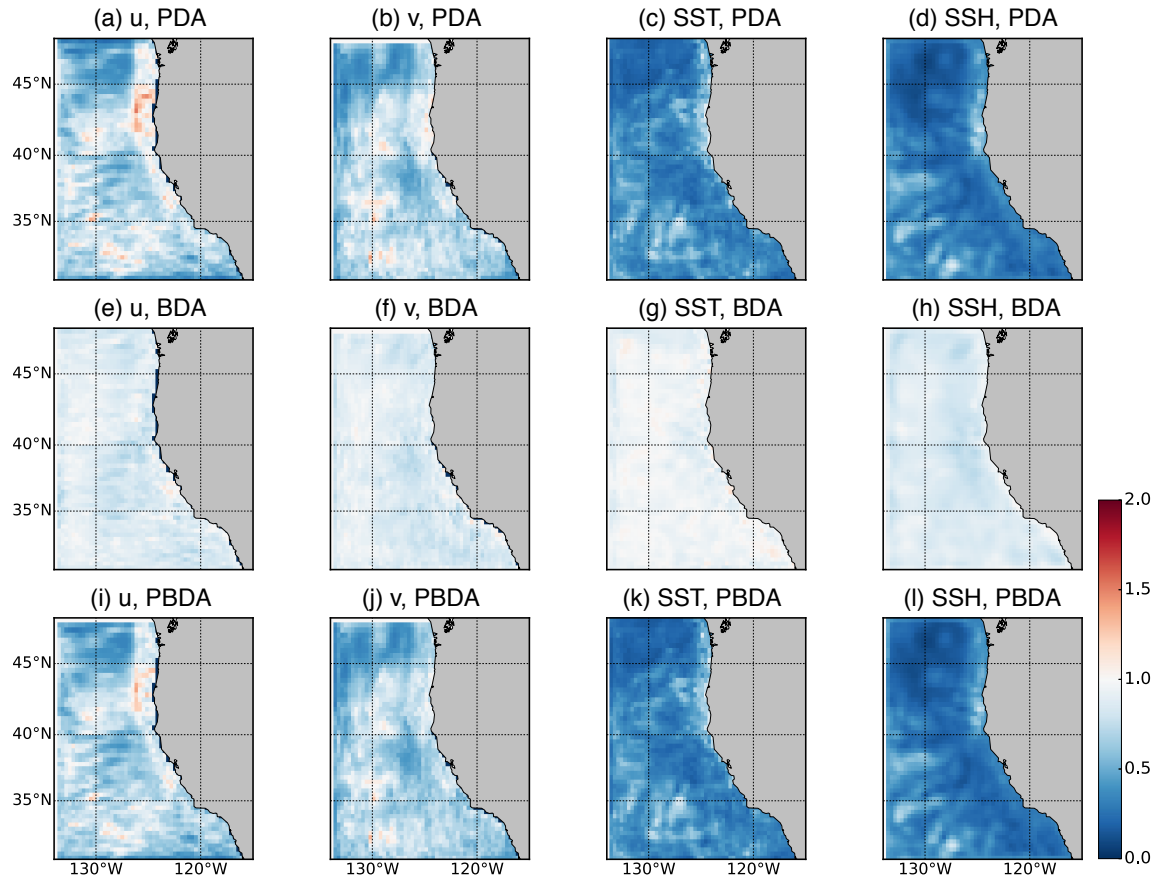


Figure 3: The ratio of the physical variables' RMSEs between data assimilation runs and free run. Smaller than 1 (cold colors) represents the reduction of the RMSE while larger than 1 (warm colors) corresponds to the increased RMSE. White areas with the value 1 mean no change in the RMSE. Top, middle and bottom rows are for PDA, BDA and PBDA, respectively, and the columns represent u, v, SST and SSH from the left to the right.

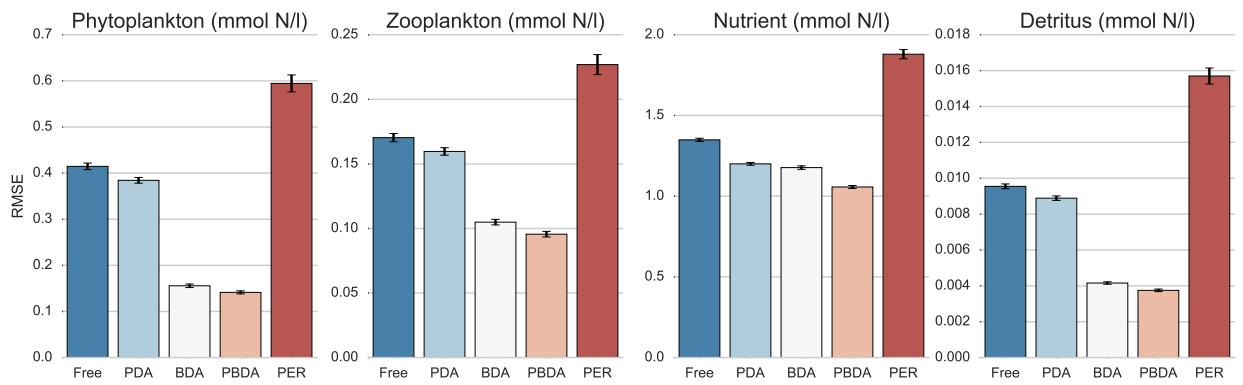


Figure 4: Same as Figure 2, but RMSE of P, Z, N and D at the surface.

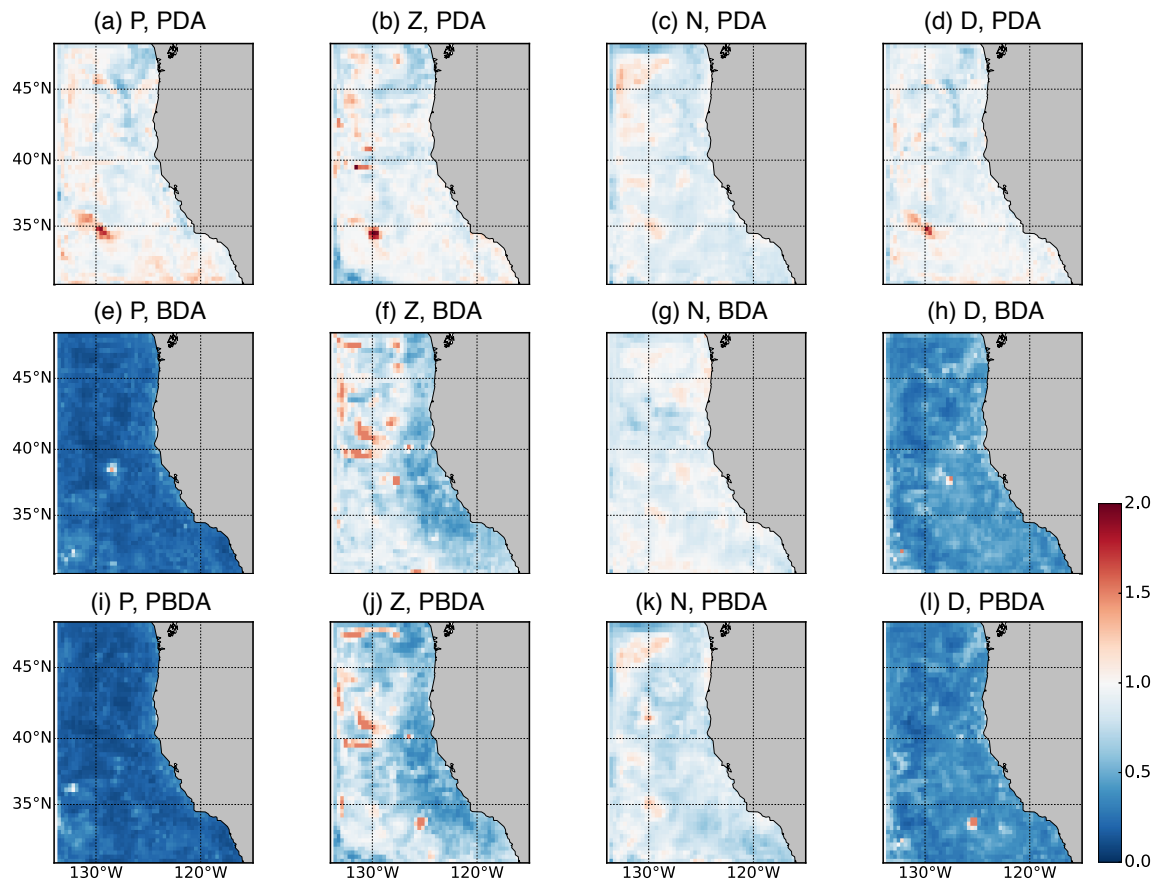


Figure 5: Same as Figure 3, but RMSE of P, Z, N and D at the surface.

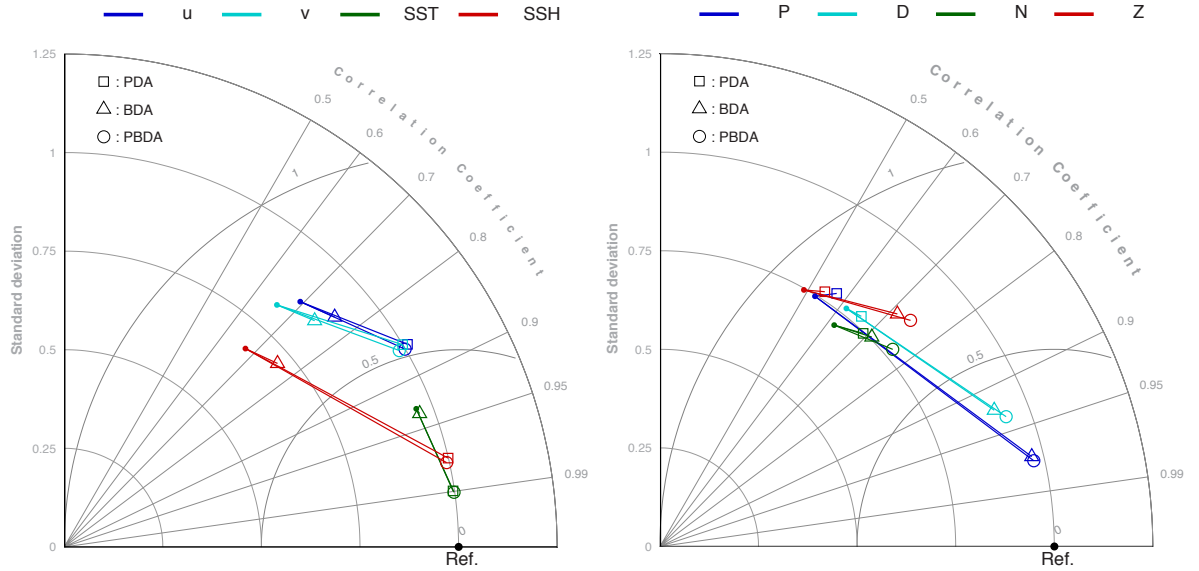


Figure 6: Taylor diagram showing the improvements of three statistical values in (left) physical and (right) biological variables. Each line shows the statistical improvement by data assimilation. The dots represent the statistical states of the prior solution, and square, triangular and circular terminators represent the statistical states of the posterior solution from PDA, BDA and PBDA, respectively. Physical variables on the right panel include u (blue), v (cyan), SST (dark green) and SSH (red). On the right panel, lines represent P (blue), Z (cyan), N (dark green) and D (red). Black dots represent the reference or true states.



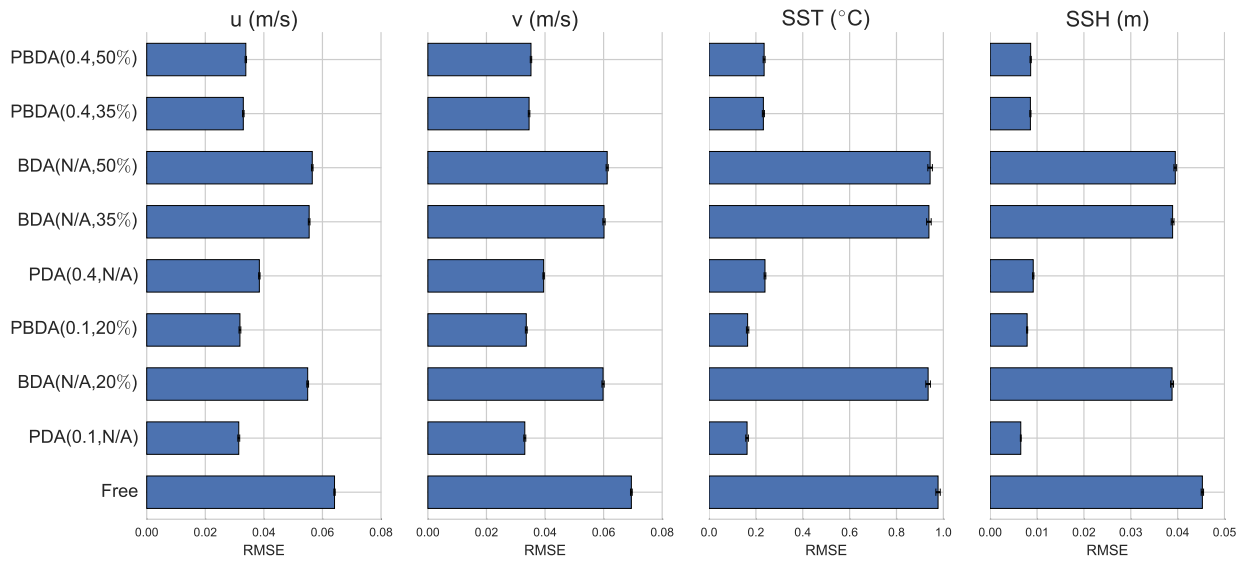


Figure 7: The RMSEs for  $u$ ,  $v$ , SST and SSH (from the left to the right) in the year 2001. The RMSEs at the bottom are from the free run, overlaid by RMSEs from data assimilation runs whose observational errors for SST and phytoplankton are indicated in the parenthesis. For example, PBDA(0.4, 50%) is the data assimilation run where the SST and phytoplankton observational errors are 0.4°C and 50%, respectively. ‘N/A’ indicates that the corresponding observation is not assimilated.

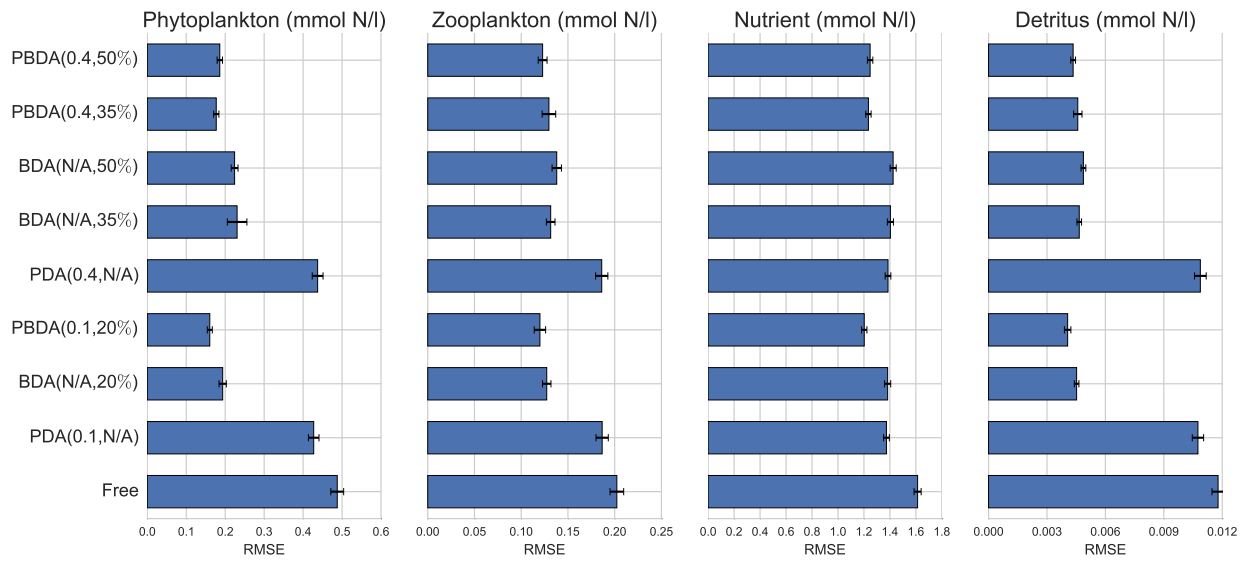


Figure 8: Same as Figure 7, but for P, Z, N and D.

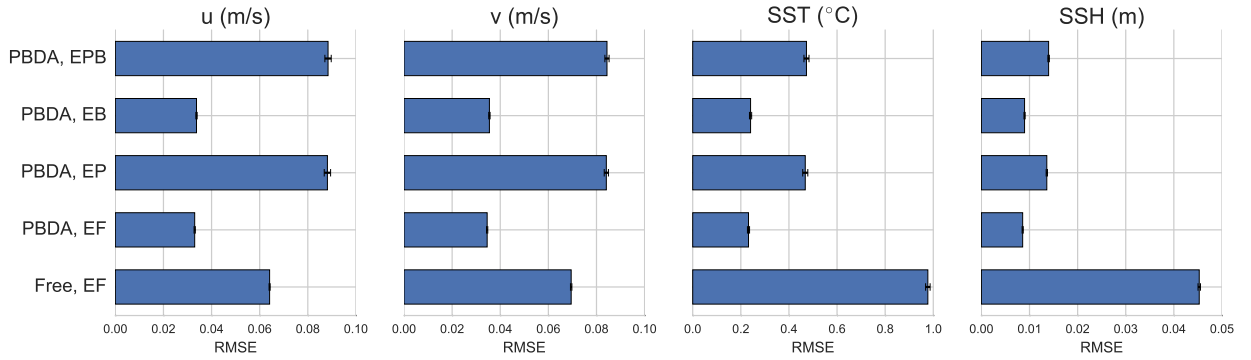


Figure 9: The RMSEs for  $u$ ,  $v$ , SST and SSH (from the left to the right) in the year 2001. The RMSEs at the bottom are from the free run with no model error, overlaid by RMSEs from data assimilation runs whose observational errors for SST, phytoplankton and the label associated with the model error are indicated in the parenthesis. The labels ‘EF’, ‘EP’ and ‘EPB’ represent ‘Error Free’, ‘Error in Physics’ and ‘Error in Physics and Biology’, respectively. The physical model error is introduced by using surface forcing of the year 2002, and the biological model error comes from different biological parameter values in the assimilation runs (Table 1).

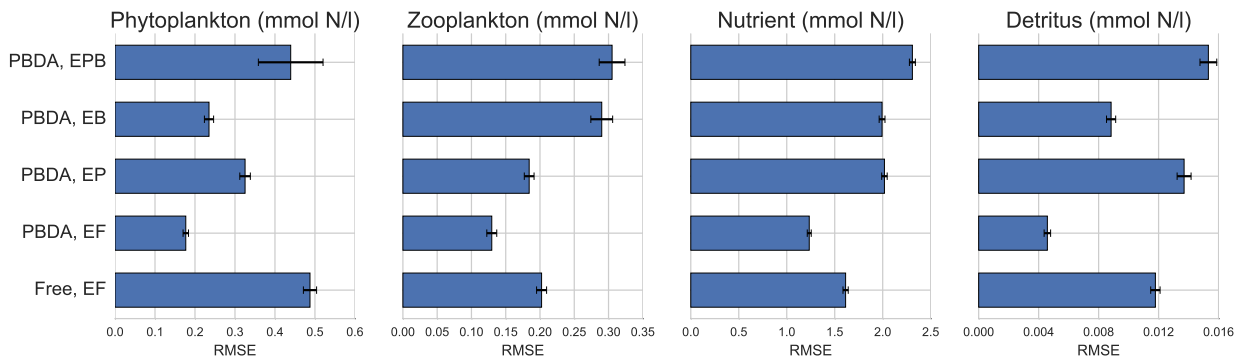


Figure 10: Same as Figure 9, but for P, Z, N and D.

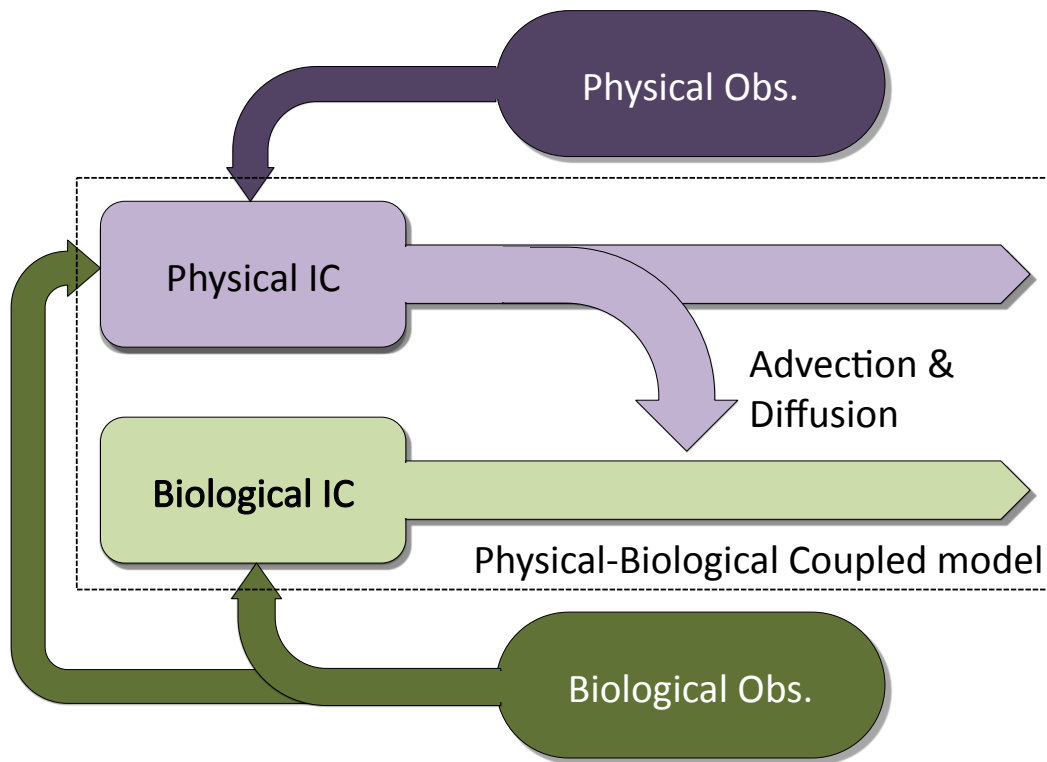


Figure 11: Diagram that shows the flow of the information from observations. Physical observations (Obs.) provide the information to adjust the physical initial condition (IC). This information is spread to biological variables through advection and diffusion. Biological observations provide the information to adjust both physical and biological initial condition. This is because the dynamics in the adjoint model pass the information only from biological component to physical component in the coupled system used in this study.