1

2 **Using Random Forests to Compare the Sensitivity of Observed Particulate Inorganic**
3 **and Particulate Organic Carbon to Environmental Conditions**
4

5 **Rui Jin[1], Anand Gnanadesikan[1], and Christopher Holder[1]**

6

7 [1] Department of Earth and Planetary Sciences, Johns Hopkins University, Baltimore, Maryland,
8 USA

9 Corresponding author: Rui Jin (ruijin@jhu.edu)

10 **Key Points:**

11 • Particulate inorganic and organic carbon (PIC and POC) estimated from satellites can be
12   predicted using environmental conditions.

13 • Random forests produce similar nonlinear relationships between some environmental
14   factors (i.e. ammonium) and PIC and POC.

15 • PIC is less sensitive to iron and more sensitive to light and mixed layer depth than POC.
16

17 **Abstract**

18 The balance between particulate inorganic carbon (PIC) and particulate organic carbon (POC)
19 holds significant importance in carbon storage within the ocean. A recent investigation delved
20 into the spatial distribution of phytoplankton and the physiological mechanisms governing their
21 growth. Employing random forests, a machine learning technique, this study unveiled apparent
22 relationships between POC and 10 environmental fields. In this work, we extend the use of
23 random forests to compare how observed PIC and POC respond to environmental conditions.
24 Our findings indicate that while both exhibit similar responses to certain environmental drivers,
25 PIC is less sensitive to iron and more sensitive to light. Intriguingly, both PIC and POC display
26 reduced sensitivity to $CO_2$, contrary to previous studies, possibly due to the elevated $pCO_2$ in our
27 dataset. This research sheds light on the underlying processes influencing carbon sequestration
28 and ocean productivity.

29 **Plain Language Summary**

30 This study looks at how different types of carbon, specifically tiny particles of chalk (particulate
31 inorganic carbon, PIC) and organic carbon from microscopic marine plants (particulate organic
32 carbon, POC), are distributed in the ocean and how they respond to environmental conditions.
33 The ratio between PIC and POC has a big impact on how carbon is stored in the ocean. We used
34 a machine learning technique to analyze how patterns in these fields estimated from satellite
35 were related to drivers such as light and nutrients. We found that PIC and POC react similarly to
36 some environmental factors (such as ammonium) but differently to others (such as iron and
37 light). Surprisingly, both types of carbon showed less sensitivity to $CO_2$ than expected from
38 previous work, possibly because of high $CO_2$ levels in the dataset.

39 **1 Introduction**

40 Because different phytoplankton functional types (PFTs) are associated with different elemental
41 cycles there is thus a need to understand how PFTs respond to different environmental drivers.
42 In particular, the ratio of particulate inorganic carbon (PIC) to particulate organic carbon (POC)
43 can play a pivotal role in the oceanic storage of carbon. POC primarily originates from
44 phytoplankton photosynthesis, resulting in the conversion of $CO_2$ into organic compounds and
45 consequent sequestration of $CO_2$ from the marine environment. Each year, nearly 10 gigatons of
46 carbon are exported from the ocean surface while around 2000 gigatons of carbon are stored in
47 the deep ocean through the biological pump (Boyd et al., 2019). However, the production of PIC
48 by calcifying planktonic organisms (e.g., coccolithophores) results in an opposing effect on
49 surface water $pCO_2$ as the accompanying reduction in seawater alkalinity leads to the release of
50 $CO_2$ (Liang et al., 2023; Kwon et al., 2009).

51 Extensive investigations have focused on deciphering the attributes of the PIC:POC ratio to
52 unravel the ramifications of global climate change on the dynamics of the oceanic carbon cycle
53 (Sarmiento et al., 2002; Rivero-Calle et al., 2015; Krumhardt et al., 2017). Archer et al. (2000)
54 argue that a decline in the PIC:POC export ratio may have contributed to the reduction in
55 atmospheric $CO_2$ that occurred during the last ice age. Brovkin et al. (2019) suggest that the
56 increase in atmospheric $CO_2$ during the Holocene was associated with changes in the rain ratio
57 and carbonate burial. Because of this, gaining a comprehensive understanding of the

58    distributional characteristics and sensitivities of PIC in comparison to POC is essential for
59    improved modeling of marine ecosystems and their responses to environmental changes.

60    In a recent investigation, Holder and Gnanadesikan (2021) utilized machine learning techniques
61    to reveal apparent relationships between the spatial distribution of phytoplankton and the
62    physiological mechanisms controlling their growth. These apparent relationships (those found in
63    the environment where many environmental drivers co-vary and where many species are present)
64    are different from intrinsic relationships found in laboratory settings where one variable at a time
65    is considered, usually for one species. Holder and Gnanadesikan (2023, henceforth HG23) found
66    that a large fraction of variability in observations can be linked to large-scale environmental
67    variables via these apparent relationships. The dominant predictors in the observational data sets
68    of POC were shortwave radiation and dissolved iron, with temperature and ammonium also
69    relatively important. However, they did not consider the impact of different physiological
70    mechanisms on different types of phytoplankton.

71    The present study juxtaposes the apparent relationships between environmental drivers of global
72    PIC and POC, allowing an assessment of how the spatiotemporal distributions of POC and PIC
73    are controlled differently. Our findings demonstrate PIC and POC exhibit distinct sensitivities to
74    variations in light, iron, and mixed layer depth.

75    **2 Methods**

76    2.1 Observations

77    A large portion of the observational data used in our analysis was compiled as part of the HG23
78    manuscript. For clarity and to minimize the requirements of the reader to seek out additional
79    scientific papers, we provide a brief overview of how the observations were compiled in HG23
80    below. For additional information on the dataset construction, please see HG23.

81    We employed observational datasets based on remote sensing as target datasets. Using remotely
82    sensed data does introduce potential sources of error into our analytical framework, as the
83    algorithms used to generate these products may be biased. However, using satellite-based
84    measurements is integral to our research objectives. First, this enables the sampling of a wide
85    range of environmental conditions while maintaining measurement consistency, thereby
86    optimizing the identification of variables that explain a substantial proportion of variance.
87    Second it facilitates the generation of datasets that are large enough for applying tree-based
88    analytical methods designed to uncover nonlinear relationships.

89    The first of these datasets was the MODIS-Aqua POC product  (Stramski, et al. 2008). This
90    particular dataset predicts POC concentrations from the remote sensing reflectances Rrs
91    measured at wavelengths of 443 and 555 nm using the equation:

92 $$POC \ = A_1[R_{rs}(443)/R_{rs}(555)]^{B_1}$$

93    Where $A_1$ and $B_1$ are regression coefficients.

94    The second target dataset utilized in our study was PIC product from Balch et al. (2005) and
95    Gordon et al. (2001). The PIC algorithm is a hybrid of two independent approaches, defined as

96  the 2-band approach and the 3-band approach. The 2-band approach uses normalized water-
97  leaving radiances in two bands near 443 and 555 nm. The 3-band approach uses spectral top-of-
98  atmosphere reflectances at three wavelengths near 670, 750, and 870 nm.

99  We accessed both PIC and POC products with a spatial resolution of 9 km and a monthly
100 climatology spanning from July 2002 to December 2022 from the NASA Ocean Color website.
101 To enhance data quality and spatial coverage, we regridded both datasets to a spatial resolution
102 of 1°.

103 2.2 Environmental drivers

104 HG23 sourced 1° monthly averaged, objectively analyzed, temperature, salinity, mixed layer
105 depth, silicate, phosphate, and nitrate from the World Ocean Atlas (WOA) 2018 dataset (Garcia
106 et al., 2019; Locarnini et al., 2019; Zweng et al., 2019). Monthly vertical velocity data at a depth
107 of 55 meters were acquired from the Estimating the Circulation and Climate of the Ocean
108 (ECCO) reanalysis dataset, version 4 release 4 (ECCO Consortium et al., 2021a, 2021b; Forget
109 et al., 2015). Net shortwave radiation (QSW) at the ocean surface from the International Satellite
110 Cloud Climatology Project (ISCCP) provided by the Objectively Analyzed Air-Sea Fluxes
111 (OAFlux) Project (Yu et al., 2006),was used as a proxy for light supply as in accordance with the
112 rationale outlined in HG23. We also use the globally interpolated MPI-ULB-SOMFFN
113 climatological $pCO_2$ product (Landschützer et al. 2020b) as an additional environmental driver.
114 No globally interpolated observational datasets are available for dissolved iron and ammonium,
115 both sparsely sampled variables. To address this, HG23 generated synthetic "observational"
116 datasets by utilizing the ensemble average of CMIP6 Earth System Models (ESMs). Both of
117 these synthetic predictors ended up being important predictors of observed POC in HG23.

118 Phytoplankton can persist under low light levels, including high-latitude areas during winter,
119 where they often enter a dormant state. Models are capable of sustaining low levels of biomass in
120 such conditions. However, the observational datasets derived from passive satellite products lack
121 information in these specific regions, resulting in an analytical gap. To address this limitation,
122 we incorporated the low-light regions into our analysis by replacing missing months at points
123 which had some measurements in the POC and PIC datasets with the 1st percentile value within
124 the corresponding global dataset (while HG23 used the 5th percentile, this difference does not
125 significantly impact the results).

126 2.3 Random Forest

127 Random Forest (RF) is a powerful ensemble learning technique widely employed in the field of
128 machine learning (Breiman, 2001). It operates by constructing a multitude of decision trees
129 during the training phase and outputs predictions based on the aggregate result of these
130 individual trees. Each tree is built on a different subset of the dataset, using a subset of
131 predictors. This contributes to its resilience against overfitting and enhances predictive accuracy.
132 Renowned for its robustness and ability to handle diverse data types, RF has become a favored
133 tool in predictive modeling, classification, and regression tasks across various domains.

134 To mitigate the risk of overfitting, we employed a random data splitting approach for both the
135 PIC and POC datasets. The dataset was split into distinct training and testing subsets with 80%

136 of the values from each dataset allocated to the training subsets and the remaining 20% forming
137 the testing subsets. This ensured that the testing subsets contained data unfamiliar to the RF
138 models during their training phase. In accordance with arguments made in HG23, decision trees
139 were constructed without sample replacement. The assessment of each RF model's performance
140 was carried out using the testing data, which were presented as input to the trained models.

141 RF models were formulated for each of the satellite-based observational estimates. The target
142 data consisted of logarithmically transformed POC or PIC variables. This transformation was
143 employed to reduce the undue influence of exceptionally large values, given the highly skewed
144 nature of both target variables. The predictor dataset, identified as "observational" for the RF
145 models, comprised observed values for sea surface temperature (SST), sea surface salinity (SSS),
146 shortwave radiation, nitrate, phosphate, silicate, $pCO_2$, reanalyzed values of upwelling velocity,
147 and model-ensemble estimates for iron and ammonium. These datasets were standardized to a
148 uniform $1°$ grid.

149 Since RFs employ a subset of variables for constructing each tree (in our case, 4 out of 11
150 predictors), it is imperative to ensure an adequate number of trees to capture the essential
151 nonlinear interactions required to model the target variable effectively. A total of 50 decision
152 trees were constructed for each RF, following the methods of HG23 who performed a
153 metanalysis to identify the optimal settings. The increase in the relative error when comparing
154 testing data and RF generated predicting data is relatively small (Table S1), suggesting the RFs
155 perform relatively well, capturing 88.7% and 83.9% of the variance in the total POC and PIC
156 datasets, respectively.

157 The assessment of variable importance within a dataset can be approached through various
158 methodologies. One of these is referred to as the permutation method. The permutation method is
159 a robust technique employed in statistical analysis and machine learning to assess the importance
160 of predictor variables in a model. In this method, a baseline is initially established by calculating
161 the model error using a trained algorithm, such as a RF. Subsequently, each predictor variable's
162 influence is evaluated by introducing randomness – the variable values are shuffled, creating a
163 modified dataset. This modified dataset is then presented to the trained model for predictions,
164 and the disparity between the error of these new predictions and the original error is computed
165 for each predictor variable. A substantial increase in the root mean squared error (RMSE) signals
166 that a particular variable holds greater importance, highlighting its significance in the predictive
167 process. Conversely, variables associated with marginal increments in error are considered less
168 influential. The permutation method thus provides valuable insights into the relative importance
169 of predictors.

170 Additionally, we conducted analyses involving the substitution of one predictor's value with its
171 observed median, while keeping the other predictor values in accordance with their observed
172 variations. This modified dataset was then presented to the RF model for analysis. A low
173 prediction in regions where the predictor variable is below the median value implies the potential
174 for this variable to suppress phytoplankton biomass.

175 Finally, in order to gain insights into the inherent relationships within each RF we conducted
176 sensitivity analyses. These analyses involved an exploration of the influence of individual
177 predictor variables. For example, when analyzing the sensitivity of iron, we adjusted its values to

178  span the minimum and maximum range observed in the observational dataset. At the same time
179  the other predictor variables were set to their median values (i.e. SW radiation was set to 176
180  W/m$^2$). This artificially constructed dataset was then supplied to the RF model to generate a
181  "median sensitivity".

## 3. Results and discussion

183  The distribution patterns of PIC and POC exhibit substantial disparities, both temporally and
184  spatially, as evident in Figure 1. In Northern Hemisphere winter, PIC concentrations (Fig. 1a)
185  demonstrate elevated levels in high-latitude regions of the Southern Hemisphere, gradually
186  diminishing as one approaches approximately 30°S latitude. Subsequently, there is an increase in
187  PIC concentrations near the equator, followed by a decline in values as latitudes increase in the
188  Northern Hemisphere. In contrast, POC concentrations (Fig 1d) exhibit their lowest values in
189  subtropical regions of both the Northern and Southern Hemispheres, with an augmentation
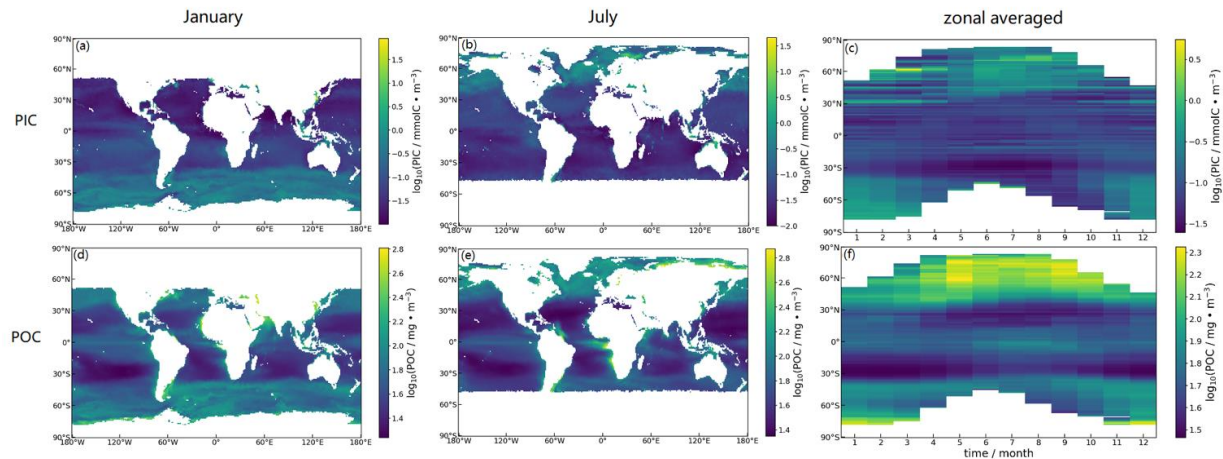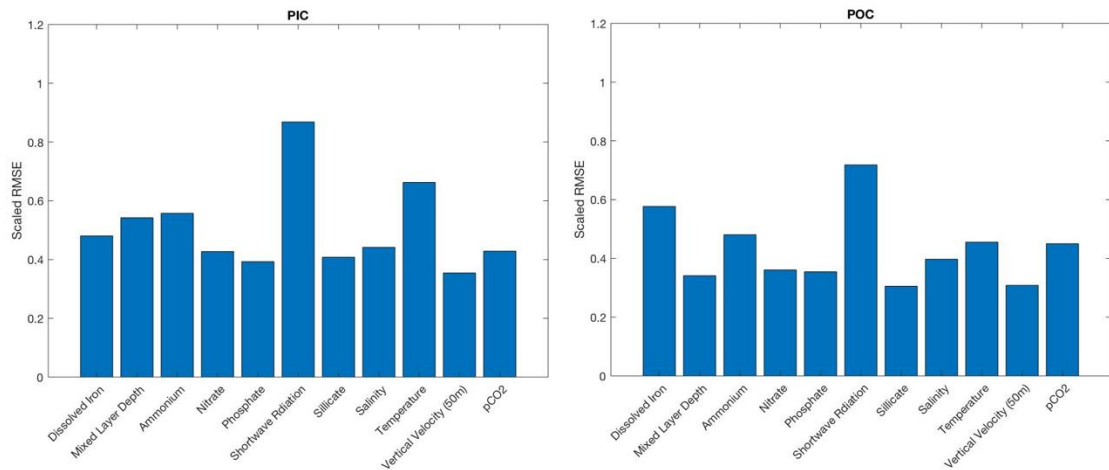190  observed around the equator and in latitudes exceeding 30°.



191

Figure 1. Global distribution of PIC and POC in January (1a and 1d) and July (1b and 1e) from
the Moderate Resolution Imaging Spectroradiometer (MODIS) averaged over all days during the
entire measuring period (2002-2022). The third column shows the zonal averaged PIC (1c) and
POC (1f). Concentrations are in log scale for better contrast.

196  Figure 1e illustrates that POC distribution in July follows a similar zonal transition pattern as
197  observed in January, albeit with different absolute values. Generally, POC concentrations in
198  high-latitude areas of the Southern Hemisphere during July are lower compared to those in
199  January, while concentrations in the Northern Hemisphere are higher. In contrast, the PIC
200  concentration in July (Fig. 1b) displays a reverse pattern when contrasted with its distribution in
201  January. During July, PIC concentrations are elevated in high-latitude regions of the Northern
202  Hemisphere, gradually declining as latitudes approach 30°S, with a minor increase near the
203  equator and reaching their lowest values in the Southern Hemisphere. Upon closer examination
204  of these distribution patterns, it becomes apparent that POC concentrations tend to align more

205  closely with the annual-mean wind stress curl field, whereas PIC concentrations are more tightly
206  coupled to seasonal changes.

207  To gain deeper insights into the contrasting distribution patterns of PIC and POC, we present
208  zonally-averaged concentrations (Fig. 1c, f). The distribution of POC concentration is
209  characterized by two distinct mid-latitude bands of reduced values, potentially attributed to
210  subsurface downwelling instigated by wind stress. Additionally, our analysis reveals that
211  between 15 and 30 degrees in both hemispheres, PIC is high during the summer and low during
212  the winter so that the peak of PIC concentration aligns with the solar zenith angle. This suggests
213  potential correlations with light, temperature or the depth of the mixed layer. It is also notable
214  that when we contrast POC and PIC in summer months for both hemispheres, a symmetry was
215  observed in PIC around 30 degrees but was not seen for POC. Near-equatorial (15°S-15°N)
216  regions show interesting differences. At 15°S, we can see a band of high values throughout the
217  year. Additionally, we see a peak that moves northward during the spring, and southward during
218  the fall, following the sun. POC shows a peak on the equator during Northern summer.

219



220      Figure 2. Variable importance plots for PIC (left) and POC (right) of the log10 transformed
221  target datasets. The x-axis shows the variables that were used in each random forest (RF). The y-
222   axis shows the relative importance of each variable computed by permuting each variable in the
223    testing dataset with the others held at their observed values, computing the root mean squared
224    error associated with the permuted inputs and normalizing this by the standard deviation of the
225                                    target from each dataset.

226  To get a better sense of the underlying determinants of PIC and POC variability, the permutation
227  importance (defined as the error when one variable is permuted for the testing data normalized
228  by the standard deviation of target data) was computed for successive variables. Large error
229  (RMSE) is indicative of predictors possessing greater importance, contributing significantly
230  towards the predictions while small error means less importance. Plots are shown in Figure 2.
231  Both datasets show that downward shortwave radiation is the most important variable. However,
232  iron is the second-most important variable in the POC data set but is only the fifth most
233  important in the PIC data set, ranking behind short wave radiation, temperature, mixed layer

234  depth and ammonium. Mixed layer depth is more important for PIC than for POC. Salinity and
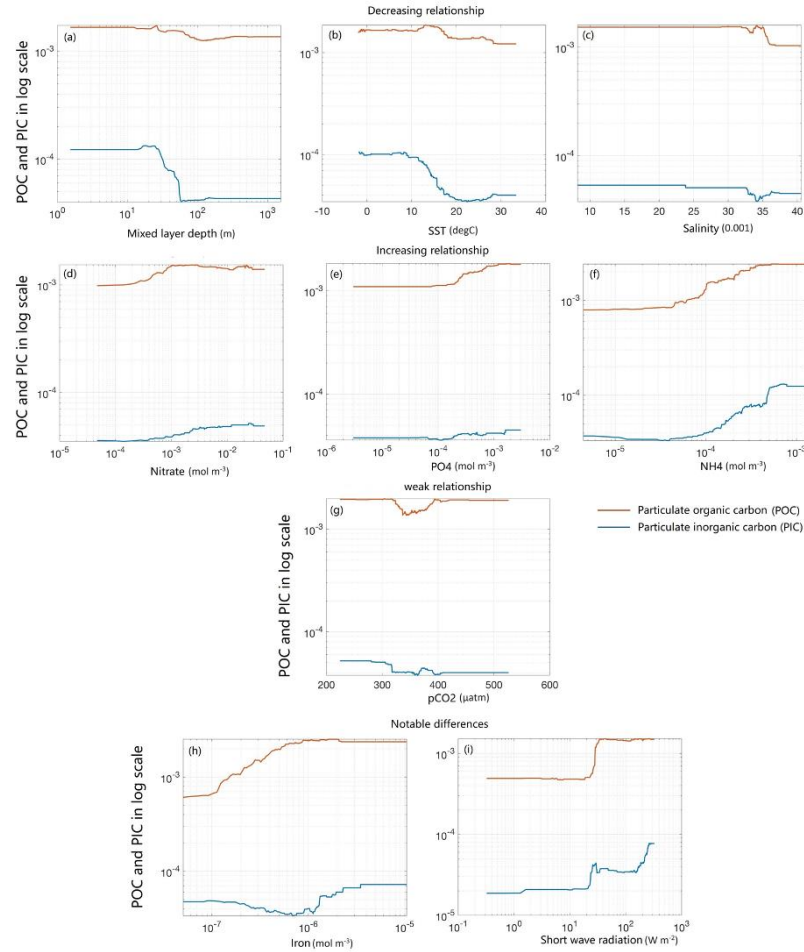235  vertical velocity are not very important in both datasets.



236

Figure 3. Sensitivity analyses on RFs trained on log10-transformed PIC (blue line) and
POC (red line) target datasets. The minimum-maximum range for each variable was
determined using values from the observational datasets and all other variables are set to
their median value.

241  We then evaluate sensitivity of PIC versus POC to individual environmental parameters with all
242  other variables fixed at their median. The first row shows that when increasing mixed layer
243  depth, temperature and salinity (Fig. 3a, 3b and 3c), both PIC and POC remain relatively stable
244  for some time then decrease at around the same concentration of the variable. For salinity,
245  however, the drop in PIC reverses when salinity concentration increases to higher values. PIC is
246  also more sensitive to changes in mixed layer depth than POC, consistent with the permutation
247  importance in Fig. 2. Conversely, greater nitrate, phosphate and ammonium (Fig.3d, 3e and 3f)
248  are associated with increases in both PIC and POC before plateauing at high values. Both PIC
249  and POC are relatively insensitive to silicate and vertical velocity as shown in the supplement
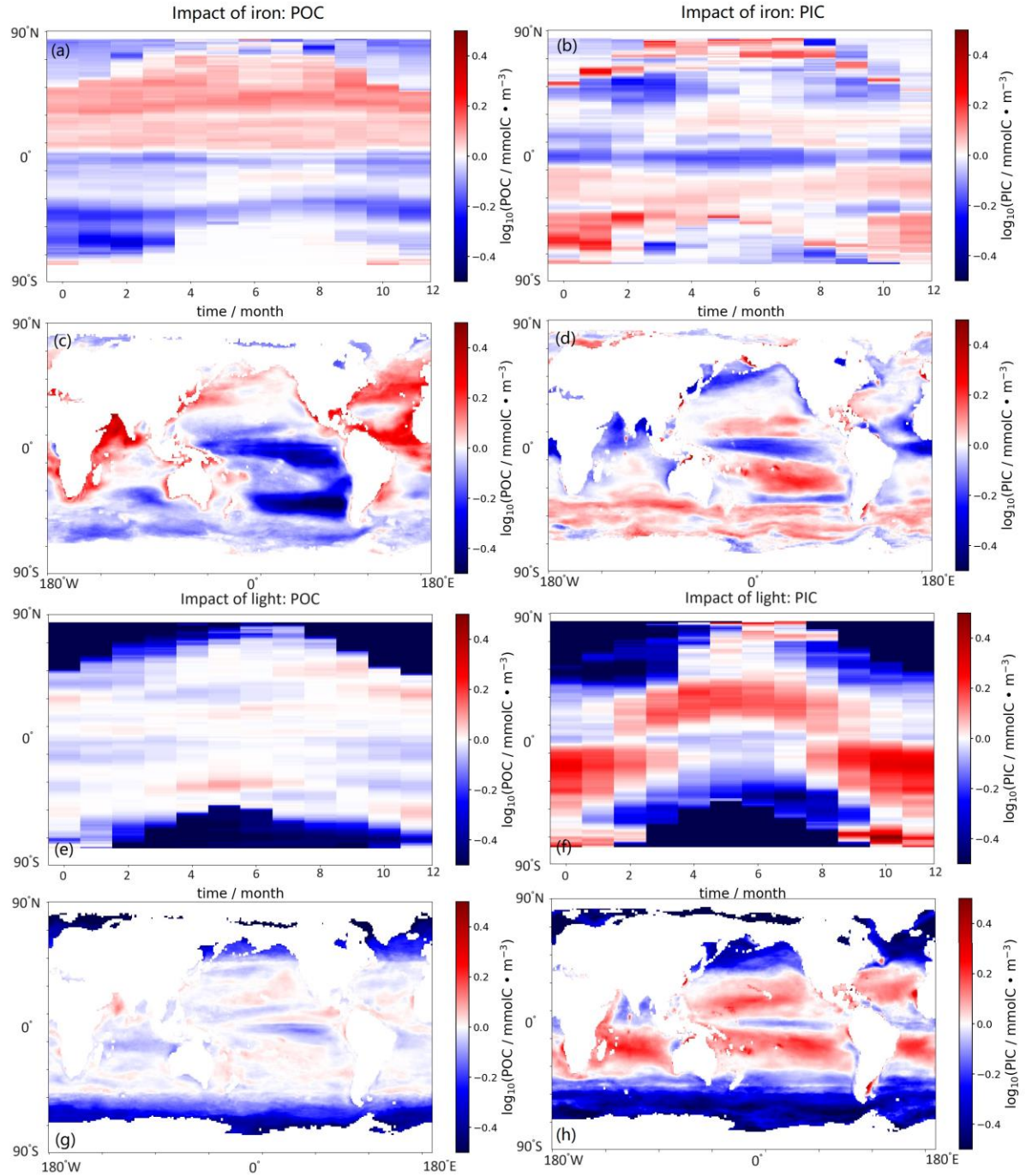
250   (Fig. S1). Both POC and PIC show relatively weak, and inconsistent, responses to changes in
251   pCO$_2$ (Fig. 3g).

252   Intriguingly, our investigation reveals distinct responses of PIC and POC to variations in iron
253   and light consistent with Fig. 2. For dissolved iron (Fig. 3h), POC shows an increase with
254   increasing iron before eventually plateauing while PIC shows a slight drop before returning to
255   the previous values. As shown in Figure 3i, POC and PIC show similar patterns between values
256   of 10 and 30 W/m$^2$, with a jump in each field observed as radiation increases. POC then reaches
257   a plateau while PIC shows a second jump around 100 W/m$^2$ as shortwave radiation increases to
258   higher values.

259   To elucidate the underlying mechanisms, we conducted a deeper examination of the spatial and
260   temporal impacts of iron and light on PIC and POC. The influence of iron on the zonally
261   averaged cycle of POC exhibits pronounced hemispheric asymmetry. In the Southern
262   Hemisphere MODIS observations, low iron levels (Fig. 4a) suppress the summertime bloom,
263   peaking in February at approximately 60°S with a 0.3 log unit reduction (roughly a factor of 2).
264   Conversely, in the Northern Hemisphere MODIS observations, higher iron levels are associated
265   with a stronger bloom, with peak enhancement occurring in May and June in subpolar latitudes,
266   also roughly a factor of two. The zonally averaged cycle of PIC under the impact of iron displays
267   different trends (Fig. 4b). Although reduced iron concentrations around the equator seem to
268   suppress PIC consistently throughout the year (consistent with POC), the results show an
269   opposite-sign sensitivity to iron compared to POC in other areas. In the Southern Hemisphere,
270   spatiotemporal iron variability fosters a more robust PIC bloom, peaking around 60°S in
271   February. In the Northern Hemisphere, higher iron levels suppress PIC around 50°S, particularly
272   in March, while iron variability promotes a PIC bloom near the Arctic region.

273   The observed annual mean impact of iron (Fig. 4c) aligns with the zonally averaged cycle,
274   revealing the most significant annual-mean biomass suppression (0.6 log units or a factor of 4) in
275   the Southeast Pacific—a region known for low iron and biomass levels (Bonnet et al., 2008), as
276   well as at the equator. Notably, higher iron emerges as a crucial factor in explaining elevated
277   POC along the boundary of the subtropical/subpolar gyre in the North Pacific, North Atlantic,
278   and the Arabian Sea. Conversely, the annual mean impact of iron on PIC has less of the ocean
279   showing strong effects. More strikingly, the spatial pattern of PIC sensitivity to iron (Fig. 4d) is
280   the opposite direction compared to POC in North Pacific, North Atlantic (particularly under the

Saharan dust plume) and Arabian Sea. Iron is associated with higher PIC levels in most parts of the Southern Ocean, as well as the South Pacific subtropical gyre.



Figure 4. Impact of variability of iron or light on POC and PIC computed by replacing the observed/modeled value at each point in time and space by the median value from observations, running the RF for each dataset, and computing the difference between the RF using the observed/modeled value and that using the observed median. Scale is $\log_{10}$, so that a value of +0.1 means that the differences between the value of iron seen at that

289  latitude and longitude and the median value of iron increases biomass by $\log_{10}(0.1)$ or
290  26% when averaged across all months.

291

292  Compared to PIC, both the zonally averaged cycle and annual mean of POC exhibit weaker
293  changes under the influence of light, with suppression observed at higher latitudes (Fig. 4e and
294  f). This observation aligns with our findings in Figure 3. The zonally averaged cycle of PIC
295  under the impact of light manifests clear hemispheric symmetry, with PIC blooms occurring in
296  both hemispheres during summer.

297  **4. Conclusions**

298  In conclusion, our study highlights divergent sensitivities of PIC and POC to distinct drivers,
299  with iron and light exhibiting particularly disparate impacts. Our findings align with the
300  conclusions summarized by Krumhardt et al. (2017) that the sensitivity of POC and PIC to iron
301  can be influenced by several factors, including temperature, $CO_2$ concentration, and the specific
302  species of coccolithophore. This opposite-sign sensitivity suggests grazing dynamics might be
303  different for PIC versus POC. In locations where sensitivity goes in the opposite direction, PIC-
304  producers and non-PIC producers might have grazers in common, so that increases in the non
305  PIC-producing phytoplankton would lead to more grazers and higher grazing pressure on the
306  PIC-producing phytoplankton.

307  Our findings show evidence for different sensitivity to light and mixed layer depth. Specifically,
308  as illustrated in Figure 3, PIC exhibits heightened sensitivity to light and mixed layer depth at
309  higher ranges, surpassing the corresponding sensitivities of POC. Furthermore, our analysis, as
310  depicted in Figure 4, demonstrates that the mean impact of light variability on PIC is notably
311  more pronounced than that on POC. These findings align with Iglesias-Rodríguez et al. (2002)'s
312  argument that a critical irradiance between 25 and 150 μmol quanta $m^{-2}\,s^{-1}$ selectively influences
313  upper ocean large-scale coccolithophorid blooms.

314  Both the PIC and POC exhibit diminished sensitivity to $CO_2$ in contrast to observational
315  syntheses made by Rivero-Calle et al. (2015) and Krumhardt et al. (2017). This divergence may
316  be attributed to our examination of comparatively contemporary data characterized by elevated
317  partial pressure of $CO_2$ (p$CO_2$) with concentrations from 325 to 407 ppmv representing the 5%-
318  95% range in our dataset. In contrast, Rivero-Calle et al. (2015) only found growth rates falling
319  when p$CO_2$ dropped below 300 ppmv, while Krumhardt et al. (2017) identified this decline at
320  concentrations below 200 ppmv. Additionally, the disparity in findings may arise from the
321  distinction in focus, with Krumhardt et al. (2017) concentrating on intrinsic relationships,
322  whereas our investigation pertains to apparent relationships.

323  Future work should aim to deepen our understanding of the intricate interplay between iron, light
324  and the dynamics of PIC and POC in marine ecosystems. Exploring the nuanced mechanisms
325  governing the response of these carbon pools to varying environmental conditions will be crucial
326  for refining predictive models and enhancing our ability to anticipate the repercussions of
327  climate change on oceanic biogeochemistry. This supports the work of Krumhardt et al. (2017)

328 who pointed out a lack of conclusive physiological responses to irradiance changes and
329 insufficient physiological data for major coccolithophore species.

330 Investigations into the specific physiological responses of key coccolithophore species to
331 fluctuations in irradiance and iron availability could provide valuable insights into the underlying
332 processes influencing carbon sequestration and ocean productivity. Long-term observational
333 studies and the integration of advanced modeling techniques may further elucidate the complex
334 relationships between environmental drivers and carbon cycling.

335

336 **Acknowledgments:** We thank Anastasia Romanou for comments on an earlier version of this

337 manscript.

338

339 **Data Availability Statement**

340 Particulate organic carbon and inorganic carbon are from the MODIS satellite climatology

341 served at NASA MODIS Climatology (NASA MODIS POC Climatology, 2020; Balch et al.,

342 2005; Gordon et al., 2001; https://oceancolor.gsfc.nasa.gov/l3/). $pCO_2$ is taken from MPI-ULB-

343 SOMFFN climatological product (Landschützer et al. 2020b). Following HG23, observations of

344 temperature, salinity, nitrate, phosphate, and silicate are taken from the World Ocean Atlas

345 (Garcia et al., 2019; Locarnini et al., 2019; Zweng et al., 2019). Shortwave radiation is taken

346 from the WHOI OAFlux data set (Yu et al., 2006). Upwelling data are taken from ECCO

347 Consortium (2021a). A compiled (climatologically averaged and aligned) data set plus a script to

348 generate the random forest and sensitivities will be available on Zenodo.

349

350 **References**

351 Archer, D., Winguth, A., Lea, D., & Mahowald, N. (2000). What caused the glacial/interglacial

352 atmospheric pCO2 cycles? Reviews of Geophysics, 38(2), 159–189.

353 https://doi.org/10.1029/1999RG000066

354    Balch, W. M., Howard R. Gordon, B. C. Bowler, D. T. Drapeau, and E. S. Booth. (2005).

355    Calcium carbonate measurements in the surface global ocean based on Moderate-Resolution

356    Imaging Spectroradiometer data. Journal of Geophysical Research: Oceans 110, no. C7.

357    https://doi.org/10.1029/2004JC002560

358    Barber, R. T., & Hiscock, M. R. (2006). A rising tide lifts all phytoplankton: Growth response of

359    other phytoplankton taxa in diatom-dominated blooms. Global Biogeochemical Cycles, 20(4).

360    https://doi.org/10.1029/2006GB002726

361    Bonnet, S., Guieu, C., Bruyant, F., Prášil, O., Van Wambeke, F., Raimbault, P., et al. (2008).

362    Nutrient limitation of primary productivity in the Southeast Pacific (BIOSOPE cruise).

363    Biogeosciences, 5(1), 215–225. https://doi.org/10.5194/bg-5-215-2008

364    Boyd, P. W., Claustre, H., Levy, M., Siegel, D. A., & Weber, T. (2019). Multi-faceted particle

365    pumps drive carbon sequestration in the ocean. Nature, 568(7752), 327–335.

366    https://doi.org/10.1038/s41586-019-1098-2

367    Breiman, L. (2001). Random Forests. Machine Learning, 45, 5–32.

368    https://doi.org/10.1023/A:1010933404324

369    Brovkin, V., Lorenz, S., Raddatz, T., Ilyina, T., Stemmler, I., Toohey, M., & Claussen, M.

370    (2019). What was the source of the atmospheric CO increase during the Holocene?

371    Biogeosciences, 16(13), 2543–2555. https://doi.org/10.5194/bg-16-2543-2019

372    ECCO Consortium, Fukumori, I., Wang, O., Fenty, I., Forget, G., Heimbach, P., & Ponte, R. M.

373    (2021a). ECCO central estimate (version 4 release 4) [Dataset]. Retrieved from

374    https://podaac.jpl.nasa.gov/dataset/ECCO_L4_OCEAN_VEL_05DEG_MONTHLY_V4R4

375   ECCO Consortium, Fukumori, I., Wang, O., Fenty, I., Forget, G., Heimbach, P., & Ponte, R. M.

376   (2021b). Synopsis of the ECCO central produc- tion global ocean and sea-ice state estimate,

377   version 4 Release 4 [Dataset]. Zenodo. https://doi.org/10.5281/zenodo.4533349

378   Forget, G., Campin, J.-M., Heimbach, P., Hill, C. N., Ponte, R. M., & Wunsch, C. (2015). ECCO

379   version 4: an integrated framework for non-linear inverse modeling and global ocean state

380   estimation. Geoscientific Model Development, 8(10), 3071–3104. https://doi.org/10.5194/gmd-

381   8-3071-2015

382   Garcia, H. E., Weathers, K. W., Paver, C. R., Smolyar, I., Boyer, T. P., Locarnini, et al. (2019).

383   World Ocean Atlas 2018, Volume 4: Dissolved inorganic nutrients (phosphate, nitrate and

384   nitrate+nitrite, silicate. Retrieved from https://www.ncei.noaa.gov/access/world-ocean-atlas-

385   2018/bin/woa18.pl

386   Gordon, Howard R., G. Chris Boynton, William M. Balch, Stephen B. Groom, Derek S.

387   Harbour, and Tim J. Smyth. (2001). Retrieval of coccolithophore calcite concentration from

388   SeaWiFS imagery. Geophysical Research Letters 28, no. 8: 1587-1590.

389   https://doi.org/10.1029/2000GL012025

390   Holder, C., & Gnanadesikan, A. (2021). Can machine learning extract the mechanisms

391   controlling phytoplankton growth from large-scale observations? – A proof-of-concept study.

392   Biogeosciences, 18(6), 1941–1970. https://doi.org/10.5194/bg-18-1941-2021

393   Holder, C., & Gnanadesikan, A. (2023). How Well do Earth System Models Capture Apparent

394   Relationships Between Phytoplankton Biomass and Environmental Variables? Global

395   Biogeochemical Cycles, 37(7). https://doi.org/10.1029/2023GB007701

396     Hopkins, J., Henson, S. A., Painter, S. C., Tyrrell, T., & Poulton, A. J. (2015). Phenological

397     characteristics of global coccolithophore blooms. Global Biogeochemical Cycles, 29(2), 239–

398     253. https://doi.org/10.1002/2014GB004919

399     Kemp, A. E. S., & Villareal, T. A. (2013). High diatom production and export in stratified waters

400     – A potential negative feedback to global warming. Progress in Oceanography, 119, 4–23.

401     https://doi.org/10.1016/j.pocean.2013.06.004

402     Klaas, C., & Archer, D. E. (2002). Association of sinking organic matter with various types of

403     mineral ballast in the deep sea: Implications for the rain ratio. Global Biogeochemical Cycles,

404     16(4). https://doi.org/10.1029/2001GB001765

405     Krumhardt, K. M., Lovenduski, N. S., Iglesias-Rodriguez, M. D., & Kleypas, J. A. (2017).

406     Coccolithophore growth and calcification in a changing ocean. Progress in Oceanography, 159,

407     276–295. https://doi.org/10.1016/j.pocean.2017.10.007

408     Kwon, E.Y., Primeau, F. and Sarmiento, J.L., 2009. The impact of remineralization depth on the

409     air–sea carbon balance. Nature Geoscience, 2(9), 630-635. https://doi.org/10.1038/ngeo612

410     Landschützer, P., Laruelle, G. G., Roobaert, A., & Regnier, P. (2020). A uniform pCO2

411     climatology combining open and coastal oceans. Earth System Science Data, 12(4), 2537–2553.

412     https://doi.org/10.5194/essd-12-2537-2020

413     Locarnini, R. A., Mishonov, A. V., Baranova, O. K., Boyer, T. P., Zweng, M. M., Garcia, et al.

414     (2019). World Ocean Atlas 2018, Volume 1:Temperature. Retrieved from

415     https://www.ncei.noaa.gov/access/world-ocean-atlas-2018/bin/woa18.pl

416     Liang, W., Han, J., Ge, Y., Zhu, W., Yang, J., & Liu, C. (2023). High-efficiency utilization of

417     biomass and seawater resources based on a distributed system with SOFC-assisted CO2 capture:

418    Feasibility analysis and optimization. Energy Conversion and Management, 296, 117675.

419    https://doi.org/10.1016/j.enconman.2023.117675

420    Iglesias-Rodríguez, M.D., Brown, C.W., Doney, S.C., Kleypas, J., Kolber, D., Kolber, Z., Hayes,

421    P.K. and Falkowski, P.G., (2002). Representing key phytoplankton functional groups in ocean

422    carbon cycle models: Coccolithophorids. Global Biogeochemical Cycles, 16(4), 47-1.

423    https://doi.org/10.1029/2001GB001454

424    Margalef, R. (1978). Life-forms of phytoplankton as survival alternatives in an unstable

425    environment. Oceanologica Acta, 1(4), 493–509.

426    Rivero-Calle, S., Gnanadesikan, A., Del Castillo, C. E., Balch, W. M., & Guikema, S. D. (2015).

427    Multidecadal increase in North Atlantic coccolithophores and the potential role of rising $CO_2$.

428    Science, 350(6267), 1533–1537. https://doi.org/10.1126/science.aaa8026

429    Sarmiento, J. L., & Gruber, N. (2002). Sinks for Anthropogenic Carbon. Physics Today, 55(8),

430    30–36. https://doi.org/10.1063/1.1510279

431    Stramski, D., Reynolds, R. A., Babin, M., Kaczmarek, S., Lewis, M. R., Röttgers, R., et al.

432    (2008). Relationships between the surface concentration of particulate organic carbon and optical

433    properties in the eastern South Pacific and eastern Atlantic Oceans. Biogeosciences, 5(1), 171–

434    201. https://doi.org/10.5194/bg-5-171-2008

435    Yu, L., Jin, X., & Weller, R. A. (2006). Objectively analyzed Air-Sea Fluxes (OAFlux) for

436    global oceans [Dataset]. Research Data Archive at the National Center for Atmospheric

437    Research, Computational and Information Systems Laboratory. https://doi.org/10.5065/0JDQ-

438    FP94

439     Zweng, M. M., Reagan, J. R., Seidov, D., Boyer, T. P., Locarnini, R. A., Garcia, et al. (2019).

440     World Ocean Atlas 2018, Volume 2: Salinity. Retrieved from

441     https://www.ncei.noaa.gov/access/world-ocean-atlas-2018/bin/woa18.pl