

Testing the skill of a species distribution model using a 21st Century virtual ecosystem.

L.R. Bardon^{1,2}, B.A. Ward², S. Dutkiewicz³, and B.B. Cael⁴

¹University of Southern California, Los Angeles, CA, USA

²University of Southampton, UK

³Massachusetts Institute of Technology, Cambridge, MA, USA

⁴National Oceanography Centre, Southampton, UK

Key Points:

- We build a correlative species distribution model to predict the global plankton biogeography of a trait-based ecosystem model
- Predictive skill varies across test cases, with functional group, and spatiotemporally, with poor end-of-century performance
- Key sources of uncertainty are traced to sampling biases in observations, and the temporal variability in target-predictor relationships

Corresponding author: Lee Bardon, lbardon@usc.edu

Abstract

[Plankton play an important role in marine food webs, in biogeochemical cycling, and in Earth’s climate; yet observations are sparse, and predictions of how they might respond to climate change vary. Correlative species distribution models (SDM’s) have been applied to predicting biogeography based on relationships to observed environmental variables. To investigate sources of uncertainty, we use a correlative SDM to predict the plankton biogeography of a 21st Century marine ecosystem model (Darwin). Darwin output is sampled to mimic historical ocean observations, and the SDM is trained using generalised additive models. We find that predictive skill varies across test cases, and between functional groups, with errors that are more attributable to spatiotemporal sampling bias than sample size. End-of-century predictions are poor, limited by changes in target-predictor relationships over time. Our findings illustrate the fundamental challenges faced by empirical models in using limited observational data to predict complex, dynamic systems.]

Plain Language Summary

[Marine plankton communities play a central role within Earth’s climate system, with important processes often divided among different ‘functional groups’. Changes in the relative abundance of these groups can therefore impact on ecosystem function. However, the oceans are vast, and samples are sparse, so global distributions are not well known. Statistical species distribution models (SDM’s) have been developed that predict global distributions based on their relationships with observed environmental variables. They appear to perform well at summarising present-day distributions, and are increasingly being used to predict ecosystem changes throughout the 21st century. But it is not guaranteed that such models remain valid over time. Rather than wait 100 years to find out, we applied a statistical SDM to a complex virtual ocean, and trained it using virtual observations that match real-world ocean samples. This allows us to jump forward to the end-of-century to test the accuracy of our predictions. The SDM performed well at qualitatively predicting ‘present day’ plankton distributions but yielded poor end-of-century predictions. Our case study emphasises both the importance of environmental variable selection, and of changes in the underlying relationships between environmental variables and plankton distributions, in terms of model validity over time.]

1 Introduction

Plankton underpin global ocean food webs and fisheries, mediate marine biogeochemical cycles, and affect climate (Fenchel, 1988; Falkowski et al., 2008; Marinov et al., 2008; Guidi et al., 2016; Hutchinson, 1961). Their global biogeography interacts with the ocean’s inventory of nutrient elements, and its capacity to sequester CO₂ (Cerneno et al., 2008; Guidi et al., 2009; Fuhrman, 2009; Falkowski et al., 1998). Understanding present and possible future biogeographic patterns of plankton communities is therefore a key component of marine microbial research. These biogeographic patterns are affected by numerous environmental factors, including supplies of nutrients and light, ambient temperature, grazing pressure, physical circulation and water column structure, and the seasonality and variability of these drivers (Tittensor et al., 2010; Rutherford et al., 1999; Graff et al., 2016). Despite substantial efforts by observational oceanographers e.g. (Lombard et al., 2019), the vastness of the global ocean and the challenges of measuring complex microscopic plankton communities makes data-limitation inevitable.

Species distribution models (SDMs) (sometimes interchangeably referred to as ecological niche models) have been widely used to predict biogeographic distributions and fundamental niche parameters in terrestrial ecosystems, and have seen a recent surge of popularity in marine ecosystem context (Flombaum et al., 2020; Righetti et al., 2019; Benedetti et al., 2021; Melo-Merino et al., 2020). While mechanistic variants exist, the

most popular implementations of SDM seek to identify the relationships between known geographic distributions of species' and sets of environmental variables. These relationships that are typically used by SDM developers to characterise biogeography in terms of where a species could, or could not, occur (Melo-Merino et al., 2020). Correlations are extracted using a variety of empirical methods, from classical statistics to bleeding-edge machine-learning (ML), or a hybridised ensemble thereof. For example, one might seek to characterise the relationships between measures of plankton concentrations (e.g. cell counts, gene markers or biomass) and simultaneously measured environmental factors (e.g. temperature, Chl-a, nutrient concentrations). The fitted model can then be used together with satellite or large synthesis database measurements to make diagnostic predictions of plankton. When the resulting SDM performs well relative to the measured datasets, predictions of species presence/absence or concentrations are then scaled globally, e.g. see (Tang & Cassar, 2019; Barton et al., 2013; Irwin et al., 2012; Agusti et al., 2019).

However, a series of assumptions and uncertainties are incorporated into correlative SDMs, many of which go unchallenged or inadequately addressed by SDM developers. While an exhaustive overview of these assumptions and uncertainties is beyond the scope of the current work (see (Wiens et al., 2009) for a thorough assessment), some are especially pertinent to marine microbial biogeography. For example, we cannot be certain that the environmental variables included in the model are a true and complete reflection of species' niche requirements', or whether some excluded or as-yet-unmeasured dimensions might better account for the observed distributions. Additionally, it is difficult to separate correlation from causation in such complex, dynamic and highly-coupled systems. Our model might highlight sea surface temperature (SST) as the primary driver of abundance; yet it remains possible that separate factors coupled to SST – perhaps underwater solar radiation penetration or nutrient supply rates – are instead more directly linked to abundance. Thus, in this scenario, and adopting the terminology of (Holder & Gnanadesikan, 2021), the relationship between SST and abundance might be described as "apparent" while the relationship between underwater solar radiation and abundance as "intrinsic". This disconnect between cause and effect can be further complicated by trade-offs in the choice of empirical model used to build the SDM, see e.g. the inverse relationship between predictive skill and interpretability in machine learning models (Carvalho et al., 2019).

There is a growing body of research that builds correlative SDMs on a variety of statistical and machine learning models, and uses them to predict global plankton biogeography from sparse observational data, both in the present day, and many decades into the future, e.g. (Righetti et al., 2019; Ibarbalz et al., 2019; Flombaum et al., 2020; Benedetti et al., 2021). Some of the results generated by such models have been highly novel and surprising, and have diverged significantly from those generated using other methodological approaches, such as trait-based mechanistic models e.g. (Ward et al., 2014; Dutkiewicz et al., 2009, 2014; Cabré et al., 2015). This is particularly true of predicting end-of-century distributions. For instance, the neural-network-derived correlative SDM developed in (Flombaum et al., 2020) predicts an increase in picophytoplankton biomass in the future subtropical oceans, in direct contrast to mechanistic ecosystem models in e.g. (Dutkiewicz et al., 2013; Marinov et al., 2010). While it is not possible to comment on which particular modelling regime best approximates the global oceans of 2100, identifying and addressing potential sources of error would be beneficial for improving accuracy and guiding interpretation.

Thus, the goal of the current work is to investigate the effects of known assumptions and uncertainties that are 'baked into' correlative SDMs, at a time when their usage is seeing an explosion of interest. To achieve this, we set up an idealised testbed to assess the predictive capabilities of an SDM built on Generalised Additive Models (GAMs) (Hastie & Tibshirani, 1986) using the output from a mechanistic global scale ecosystem

model, the ‘Darwin’ model (Dutkiewicz et al., 2021), as a ‘ground truth’. To explore the effect of spatiotemporal biases in real-world observational datasets, Darwin model outputs are sampled in space and time to mimic historical ocean measurements, and also randomly. The resulting SDM is then evaluated in its ability to capture the virtual ocean’s emergent biogeography in the present day ‘*spatial predictions*’ and by the end-of-century ‘*temporal predictions*’. Our experiment is thus designed to generate insights into the fundamental limitations of correlative SDMs, applied in the current context, as a function of their core assumptions and uncertainties.

At the outset, we stress that our intention here is not to raise a false dichotomy whereby one particular methodological approach is pitted against another to decide a ‘winner’. Nor are we making any claim as to the accuracy of the Darwin model in its ability to faithfully predict plankton abundance and diversity in the real ocean. Rather, the following case study is designed to assess how a correlative SDM might fare in predicting a complex but well-understood microbial ecosystem (see e.g. (Dutkiewicz et al., 2020)) embedded in a dynamic, self-consistent model of the Earth’s ocean through time.

2 Materials & Methods

We performed a suite of tests using a widely applied implementation of GAMs (Servén & Brummitt, 2018) as our SDM and the Darwin model, a dynamic marine microbial ecosystem model coupled to an Earth system model ((Dutkiewicz et al., 2021), (Sokolov, 2005)). Our decision to use GAMs as the empirical framework underlying our correlative SDM was informed by the work of (Righetti et al., 2019), who demonstrated that GAMs perform comparably to Random Forest and Generalised Linear Models in a range of relevant predictive tasks, while offering a higher degree of both interpretability and flexibility. Additionally, GAMs are of intermediate complexity between classical statistical regression models, and more sophisticated machine learning methods, which arguably makes them both accessible and potentially attractive to a wide range of researchers. Nonetheless, we note that we could have selected any one of a wide variety of statistical or machine learning algorithms, each with their own unique pros and cons.

To train the GAMs, we sample the Darwin model at the same places and times as in a large ocean measurement dataset used for similar purposes (Martiny & Flombaum, 2020). The resulting GAMs SDM is then used to predict Darwin model plankton biogeography. To quantify how spatiotemporal bias in the training dataset affects predictive skill, we train an additional set of GAMs using a dataset of the same size, but sampled uniformly randomly across the virtual ocean’s surface, and uniformly randomly over the same period of time. To quantify the effect of training set sample size on predictive skill, we generate 54 additional random-sample training sets, in 18 different sample sizes. We evaluate the ability of the SDM to predict the global biogeography of the different plankton functional groups in the simulation, both during the 22-year period over which measurements were taken (i.e. spatial extrapolation), and during the last 22 years of the 21st century (i.e. both spatial and temporal extrapolation).

2.1 Numerical Model Simulation

The Darwin model ecosystem used here includes 51 plankton populations across 7 functional groups (2 prokaryotes (pro), 2 pico-eukaryotes (pico), 5 coccolithophores (cocco), 5 diazotrophs (diazot), 11 diatoms (diatom), 10 mixotrophic dinoflagellates (dino) and 16 zooplankton (zoo)). Individual populations correspond to different size classes within functional groups, with all size classes covering a range of 0.6–2425 μm equivalent spherical diameter. Functional groups have distinct allometric relationships for growth, grazing, and sinking parameters (see (Dutkiewicz et al., 2020)). The model ecosystem is embedded within the Massachusetts Institute of Technology Integrated Global System Model (IGSM) (Prinn, 2013; Sokolov, 2005) which includes modules for the physics, chemistry,

and biogeochemistry of the atmosphere, land and ocean. The ocean component has a $2^\circ \times 2.5^\circ$ resolution grid and 22 vertical layers (10m thickness at surface to 500m at bottom). The simulation is forced with observed greenhouse gas emissions from 1860–1990 and then with a high emissions scenario that is analogous to the IPCC’s Representative Concentration Pathway 8.5, from 1990 – 2110. This perturbation results in $\sim 3^\circ\text{C}$ sea surface temperature warming by 2100, sea ice retreat, increased stratification, and an altered overturning circulation. The IGSM has been used to examine changes in marine biogeochemistry and ecology in previous studies (e.g. (Dutkiewicz et al., 2013) but with a simpler version of the ecosystem model. The current more complex ecosystem has also been used in previous studies, but only for the present day’s ocean (Dutkiewicz et al., 2021; Sonnewald et al., 2020; Kuhn et al., 2019). This model and previous model validation for the present day demonstrates that the output compares well with observations along both axes of size and functional type (e.g. (Dutkiewicz et al., 2021, 2020)).

2.2 Ecosystem and Environmental Variables

Surface-level plankton abundance data and environmental parameters were extracted from Darwin simulation outputs, where surface in this context refers to the 10m thick surface grid box. The ecosystem data contains 51 separate plankton biomasses, arranged into seven functional groups (as described above). A number of environmental variables have frequently been integrated into correlative SDMs to predict abundance and diversity, and have thus been included here. They are: sea surface temperature (SST), photosynthetically active radiation (PAR), phosphate (PO_4), nitrate (NO_3), silicate (Si) and iron (Fe). We sampled both the plankton abundance data and the environmental predictor variables from the 3586 spatiotemporal cells that encompass the representative ocean measurement coordinates, and from the 3586 randomly selected spatiotemporal cells. Note that the model simulation used for the current analysis nominally starts in 1991 and extends to 2100. As such, we sample the model output from the beginning of 1991 to the end of 2012 and consider this as a substitute to 1987–2008 in this context. This is justified because the Darwin model’s internal variability does not match real-world interannual variability in terms of timing, though does capture the magnitudes (e.g. there are El Niño events, but these do not occur in the same years as the real ocean). To validate predictions, we also consider whole-ocean surface data over the same period, and for the final 22 years of the simulation, from 2079 – 2100.

2.3 Building the Correlative SDM

Although GAMs have considerable flexibility in how their core components are selected, we used the standard ‘LinearGAM’ model of the freely available PyGAM package (Servén & Brummitt, 2018). LinearGAM incorporates a Gaussian distribution function with an identity link function, and fits predictor functions using penalised B-splines. In combination, these components impose smoothness to prevent over-fitting, and enable the automatic fitting of nonlinear relationships. For an initial set of results, we set the number of permitted splines to 20 for each predictor variable. We note that our results are not sensitive to the choice of this parameter (see ‘Model Comparison & Sensitivity Tests’). At the outset, we attempted to resolve and make predictions for individual plankton tracers, but the resulting models proved to be highly unstable, so we instead choose to proceed by summing the abundance data for each functional group, and training GAMs accordingly. The resulting partial dependency plots were examined for unexpected behaviours, or any clear indications of over or under-fitting. The resulting GAMs SDM was then used to make predictions for the global surface ocean plankton biomasses during 1987–2008 and 2079–2100.

2.4 Model Comparison & Sensitivity Tests

We define presence/absence as modelled biomass being above/below a cutoff threshold (10^{-5} mmol C/m³), but find that patterns in the resulting predictions are not sensitive to the choice of this threshold (Table S4).

The R^2 value of the GAMs predictions against the ‘ground-truth’ simulation values is given as $R^2 = 1 - SS_{res}/SS_{tot}$, where SS_{res} is the residual sum of squares and SS_{tot} is the total sum of squares. While R^2 is a widely-used statistic in regression analyses, it does not by itself provide a complete picture of goodness of fit. We therefore also examine the mean and median relative differences, defined here as $\bar{X}_{me} = (mean_{predicted} - mean_{actual})/mean_{actual}$ and $\bar{X}_{md} = (median_{predicted} - median_{actual})/median_{actual}$, as an indicator of bias. We also consider the false positive and false negative fractions, i.e. the fraction of grid cells where the GAMs incorrectly predict, respectively, present and absent biomass. Finally, we performed the above analyses with the logarithm of biomass concentrations and found that our results were not sensitive to this choice. Overall, we found that coccolithophores yielded the median performance in terms of goodness of fit with respect to spatial extrapolations. As such, this group is featured in the main body of this work, while results for the other six functional groups are reported in the supplements.

GAM sensitivity was investigated by varying the number of splines used in performing the fits; first by halving to 10, and then doubling to 40. While the resulting partial dependency plots revealed a clear change to the smoothness of the fit, as expected, we found that the resulting statistics were not appreciably impacted. To investigate the effect of sample size on the overall predictive power of the GAMs, we vary the number of randomly-sampled cells from a minimum of 100 (reducing to 63 ocean cells), to a maximum of 20,000 (reducing to 11,557 ocean cells), using 18 different test cases. Each sample size test case consists of three independent random samples, with the mean value being reported along with the standard deviation (Figure 4).

We also performed a range of simpler correlation analyses, to build a broader picture of the emergent relationships between functional group biomass and predictors. These act as a visual aid to better understand how these relationships might change in time and space, and as a basic cross-reference for GAMs-derived partial dependence plots of the training sets. We first calculate the Pearson’s Correlation Coefficient (ρ) for each functional group-predictor pair, and the Spearman’s Rank Correlation Coefficient (ρ_s). Respectively, these popular methods detect the strength of linear associations between variables, and the strength of correlation in monotonic relationships. A commonly used method for addressing skew or capturing scaling relationships is the log-transform, which we apply to all datasets before recalculating ρ . However, this method of broadly applying a single transformation is not optimal. A more robust approach would be to examine the distribution of each target-predictor relationship individually, before an appropriate transformation is selected. Nonetheless, even this more optimal method runs the risk propagating transformation uncertainty into the resulting confidence interval.

With these limitations in mind, we also determine correlations using the more recent distance correlations method of (Székely et al., 2007). This technique captures the strength of both linear and nonlinear associations and avoids the need to make assumptions about variable distributions or linearity. We plot the correlation matrices for the main 3586 cell test cases, both measurements-derived and randomly-sampled, in 1987-2008, and at the same locations in 2079-2100. We explore the effect of sample size on the derived correlations by increasing the number of randomly-sampled cells to 12,894, and finally to 25,683 cells.

3 Results

3.1 Spatial Predictions

We first describe the results of predicting plankton biogeography during the historical measurement period (1987 – 2008) (Figure 1). We find that predictive ability varies considerably across functional groups. There are fewer instances of our SDM incorrectly predicting presence (false positive) or absence (false negative) biomass for prokaryotes, picophytoplankton and coccolithophores (16–19% of all location-month pairs) than for diatoms, diazotrophs, and dinoflagellates (26–31%), with zooplankton in between (21%). Where biomass is present and is predicted as such, the SDM’s predictive ability for biomass concentration also varies substantially between functional groups (Figure 2); the SDM accounts for as much as 71% of the variance in biomass (diazotrophs) and as little as 41% (zooplankton). These patterns are reflected also in the mean relative differences and the balanced accuracy.

Patterns of overprediction of biomass occurs across most of the oceans. For prokaryotes, picoeukaryotes, dinoflagellates and zooplankton, this is especially evident in the Arctic (see Figures (c) of S1, S2, S5, S6). For these groups, we also see consistent underprediction in most of the Indian Ocean and in the Eastern Equatorial Pacific. Meanwhile, diatoms are substantially overpredicted in most of the mid- and high-latitudes in the Northern Hemisphere but perform relatively well in the subtropics (Figure S4(c)). Diazotrophs yield the best overall performance, with only a small amount of overprediction in the subtropical Atlantic, and overprediction in the transition zone latitudes poleward of the subtropics (Figure S3(c)).

In general the SDM shows a tendency to overestimate biomass in the spatial predictions regime. Overestimation ranges between 9–21% on average (picoeukaryotes and zooplankton, respectively), with a median overprediction of $\geq 16\%$. Despite this, there are some notable instances in the current context where the model performs well. Spatial predictions for coccolithophores, prokaryotes and diazotrophs all yield R^2 values that range between 0.62 and 0.71 (Figures 1(e), S1(e), S5(e)). Diazotrophs fare particularly well in this regime, with a mean overprediction of 10%, an R^2 of 0.71, and the best visual, qualitative match of biogeography overall (although we note that the median overprediction in this case is a substantial 194%) (Figures S3(c) and S3(e)). On the whole, the SDM trained on data from historical measurement locations appear to be able to reproduce qualitative biogeographic patterns from spatial predictions well, but quantitative performance is variable, with a broad tendency towards overprediction. Notably, the greatest predictive errors more often occur in the undersampled regions of the ocean, such as the Arctic and Indian Oceans, but are by no means confined to these regions. For instance in the highly sampled North Atlantic predictions for diatoms and diazotrophs was also poor.

3.2 Temporal Predictions

The SDM’s predictive ability is substantially reduced when extrapolating to the future ocean (see Figures 1 and 2). Rates of false positives and negatives in presence/absence do not uniformly change across functional groups: the cosmopolitan groups whose ranges expand poleward experience the least overall change, increasing by between 3% and 11% in prokaryotes, dinoflagellates and coccolithophores, with a decrease of 5% for picophytoplankton. The SDM’s ability to correctly predict presence/absence is further reduced for the groups with a more confined biogeography, increasing by between 14% and 23% for diazotrophs, zooplankton and diatoms. We see a substantial increase in false negative occurrences for diatoms (to 29%), the group whose biogeographic range contracts most. Where biomass is present and is predicted as such, the SDM’s predictive ability was reduced for all functional groups. In most cases, this reduction is substantial, with the fraction of variance accounted for by the SDM reducing by between 17 and 50%, such

that the prediction for zooplankton is worse than just assuming a globally uniform constant biomass (i.e. $R^2 < 0$). We see a marked increase in mean relative differences compared to the ‘spatial’ predictions, accompanied by a reduction in balanced accuracy for all groups besides diatoms (Figure 2).

Diatoms are the only group for which the fraction of variance accounted for does not decrease substantially, only from $R^2 = 0.59$ to $R^2 = 0.56$ (Figure S4). Thus, the predictive ability for diatom biomass where it is present is not greatly reduced, despite the SDM’s substantial overprediction of the contraction of diatoms’ biogeography. This is not sensitive to varying the absence/presence cut-off value by an order or magnitude in either direction (Table S1).

Spatial patterns of prediction errors of coccolithophores, prokaryotes, picoeukaryotes, dinoflagellates and zooplankton are largely similar to those for the historical period, except the North Atlantic is now underpredicted for all groups besides diazotrophs (Figures 1, S1, S2, S4, S5, S6). Diatom biomass is notably underpredicted in the Southern Ocean and Northern Atlantic (Figure S4). Meanwhile, diazotroph biomass is notably overpredicted throughout the Atlantic Ocean, the Arctic, bands of the subtropical Pacific and Indian Ocean (Figure S3). Excluding diatoms, the overall tendency towards overprediction is exacerbated for all groups, increasing by 57% for prokaryotes, picoeukaryotes, coccolithophores, and dinoflagellates, by 20% for zooplankton, and by 49% for diazotrophs. Median overpredictions also increase for all groups besides diatoms.

3.3 Model Trained on Randomised Locations

Here we compared the above results with those produced when the GAMs SDM was trained on randomly sampled datasets (Figure 2). Interestingly, the broad spatial patterns of where overprediction and underprediction occurs do not change much when training the SDM on randomly distributed data, as opposed to the ocean observation locations (Figures S8 and S9). Nonetheless, predictive abilities increase, biases are reduced, and balanced accuracy increases in both the spatial and temporal cases (Figure 2). The fraction of variance accounted for by the SDM increases by 2–19% when using random data to predict historical biogeography, but increase from 5–46% when using random data to predict future biogeography. The most notable differences are for prokaryotic, picoeukaryotic, and zooplankton biomass in the future case. The magnitude of the biases also decreases – average biases are within 3–4% in the historical case using random data. The median bias for all groups is still that of overprediction, with most groups in the range of $\geq 17\%$ compared to $\geq 30\%$ for measurements-derived predictions. Diatoms and diazotrophs have a markedly higher bias in both measurements-derived and random cases, of $\geq 194\%$ and $\geq 162\%$, and $\geq 65\%$ and $\geq 35\%$. In the future case, using random data reduces biases for all groups, though does not eliminate them. We also found that the predictive ability of the SDM was only weakly dependent on sample size (where sample size here refers to the number of grid cell-month pairs that are sampled) (Figure 4), with predictive ability appearing to plateau with increasing sample size.

The results using random training datasets suggest that historical measurement biases reduce the predictive ability of the SDM more than the sample size of the training dataset. Predictive ability can be improved by subsampling or weighting one’s training dataset to reduce biases in space and time, although the coarse resolution of the Darwin model – and thus reduced variability as a result of correlated observations – relative to the real ocean may contribute to this plateauing effect.

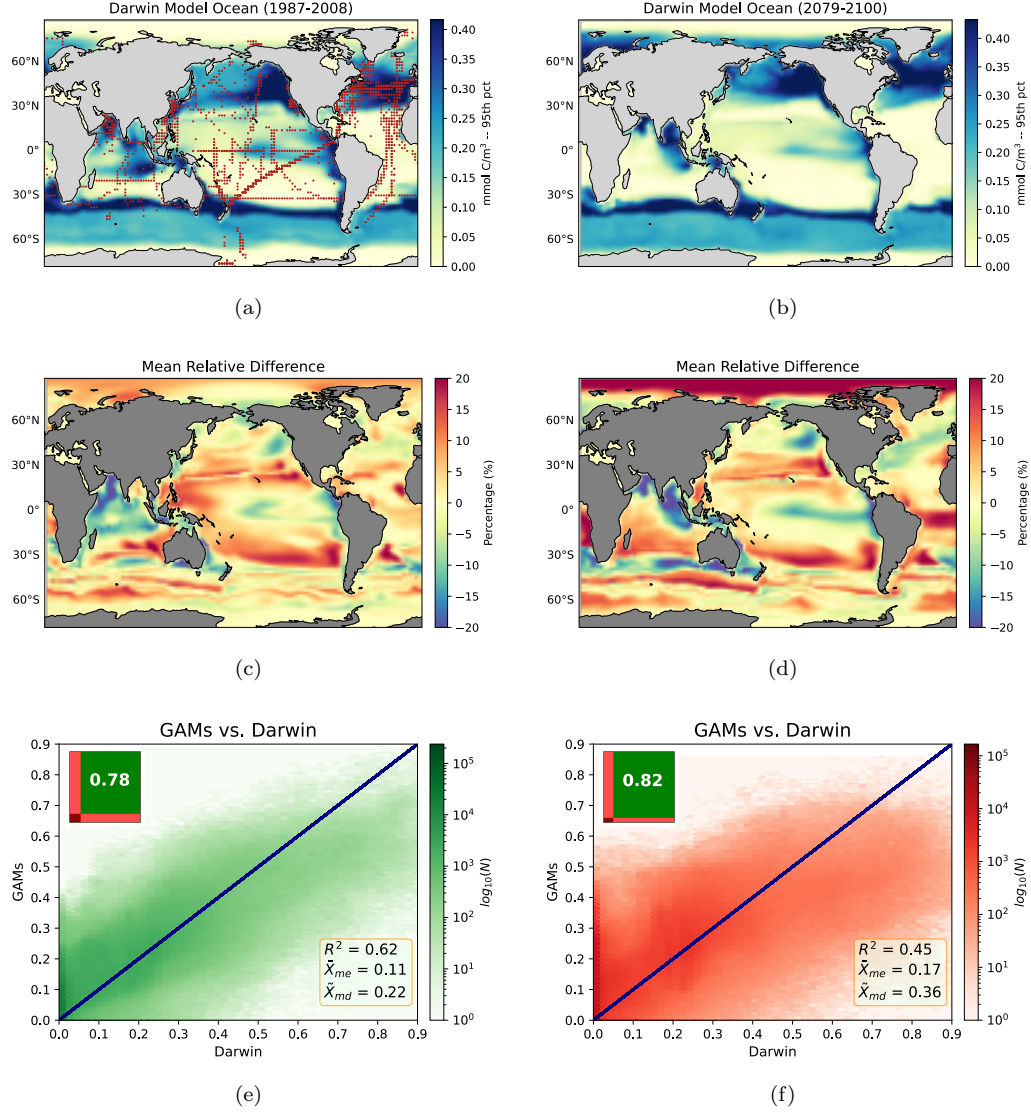


Figure 1: **(a)** Mean coccolithophore surface biomass (1987 - 2008) from the Darwin model. Red points indicate spatial location of training set datapoints, derived from ocean measurement data. **(b)** As per 1(a) for the years 2079 - 2100. **(c)** Relative (percent) difference between mean coccolithophore surface biomass from the Darwin model and the GAMs SDM (1987 - 2008) **(d)** As per 1(c) for the years 2079 - 2100. For direct visual comparison, we first calculate the 5th and 95th percentile of the relative difference values for both the spatial and temporal predictions, then scale symmetrically to whichever of these values is the greatest, in either direction. **(e)** Hexagonally binned scatterplot of 1987-2008 GAMs SDM predictions vs 1987-2008 Darwin model. Colorbar shows log-scaled density of observations. *Top inset:* Fraction of data above the presence/absence threshold (10^{-5} mmol C/m³)(green box), GAMs SDM below threshold (left, light red), Darwin below threshold (bottom, light red), both below threshold (dark red). *Bottom inset:* The R^2 , relative difference of the means (\bar{X}_{me}), and relative difference of the medians (\bar{X}_{md}). **(f)** As per 1(e) but for 2079-2100. See Supplemental Materials for other functional groups.

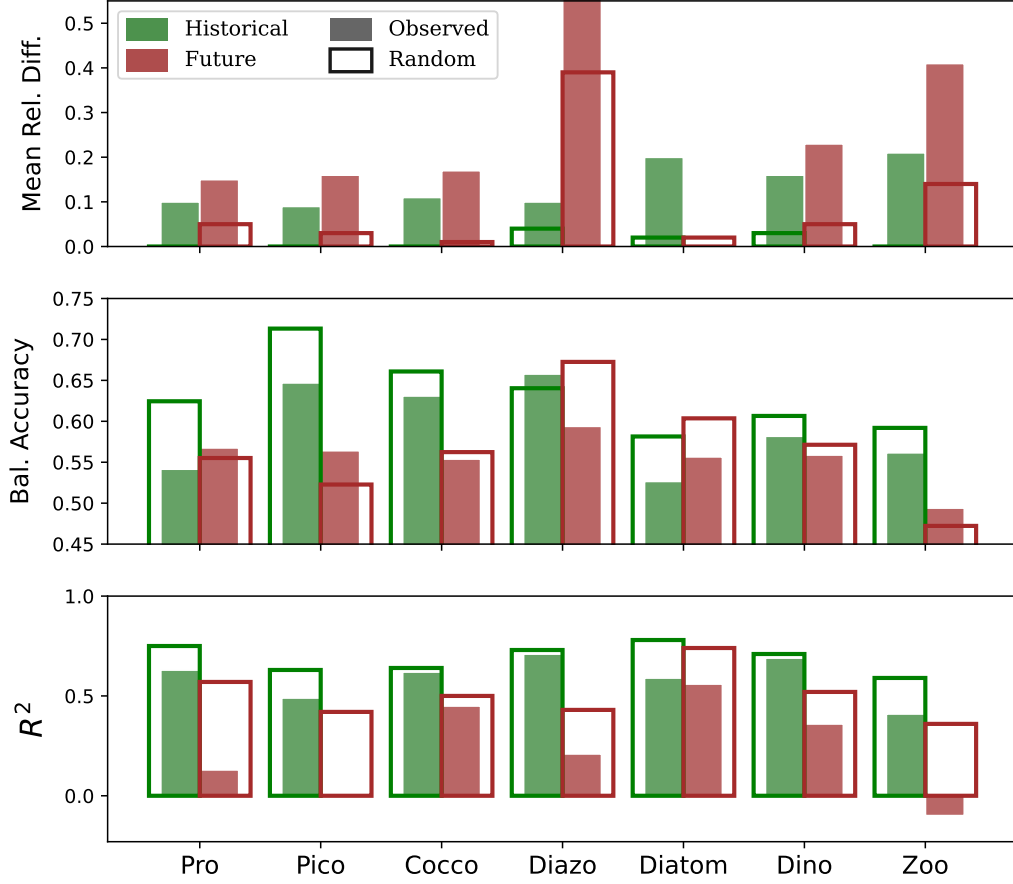


Figure 2: Comparing Darwin model ‘true’ biomasses with GAMs SDM predictions for each functional group in 1987-2008 (historical) and 2079-2100 (future), and from measurements-derived and randomly-sampled training sets. *Top to Bottom:* (a) Relative differences of the means, given by $(GAMs_{mean} - Darwin_{mean})/Darwin_{mean}$. (b) Balanced accuracy, given by $(sensitivity + specificity)/2$. (c) R^2

4 Discussion

Broadly, our GAMs-driven correlative SDM demonstrates capability in qualitatively capturing large-scale spatial patterns of plankton biogeography, but struggles to make robust quantitative predictions. This is particularly evident when the model is trained on historical ocean measurement data, and used to predict future plankton biogeography as a response to climate change. The emergent relationships between predictor variables and plankton abundances change spatially, seasonally and over the longer term. This is demonstrated by the variable nature of the partial dependence plots (Figure 3(a)–(b) and Figures S10 and S11), and by the change in correlation strengths identified by each of the independent methods used in generating the correlation matrices (Figure 3(c)–(f) and Figure S12). The correlation matrices offer an especially powerful visual demonstration of these points; we clearly see the change in apparent relationships between biomass and environmental predictors in the measurements-derived sample space, assessed over the same period of time one hundred years into the future (Figure 3(c) and 3(d)). It’s

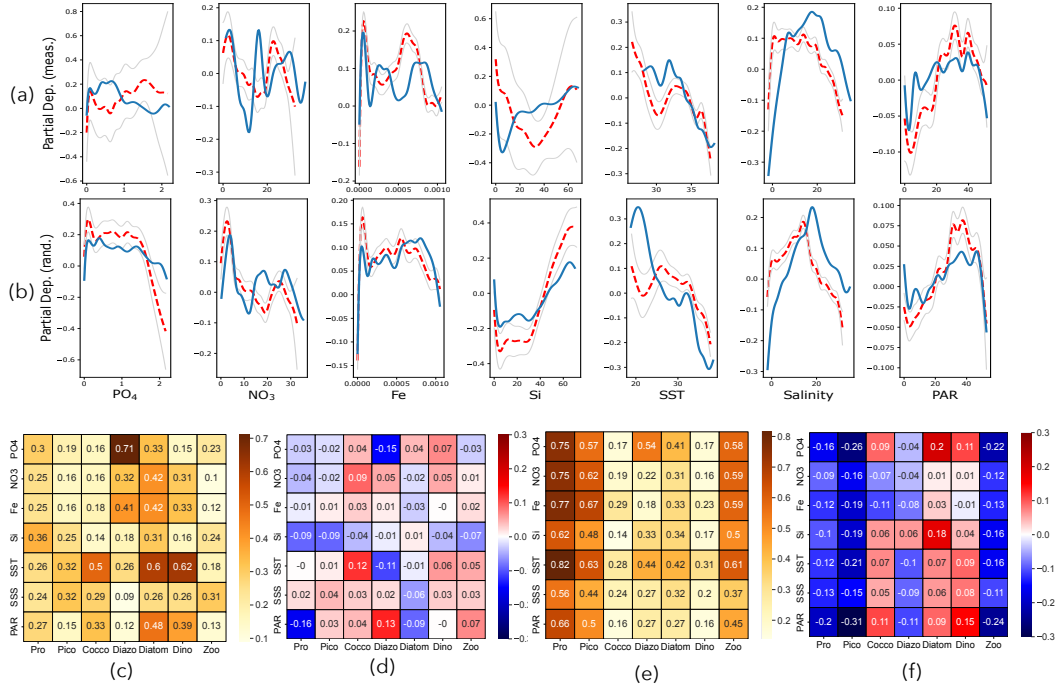


Figure 3: *Changing Relationships*: (a) Partial dependence plots of coccolithophore biomass (mmol C/m^3) as a function of each predictor, centred around the median (PO_4 , NO_3 , Fe, Si in mmol X/m^3 , SST in $^\circ\text{C}$, SSS in PSU, PAR in $\text{E/m}^2/\text{day}$). Plotted using data from 3586 Darwin surface ocean cells at measurements-derived locations spanning 1987-2008 (dashed red line) and at the same locations from 2079-2100 (blue line). Grey lines indicate 95% confidence interval for the 1987-2008 case. (b) As per 3(a), but using data from 3586 randomly sampled cells. (c) Correlation heatmap for the measurements-derived training set, 1987-2008, generated using the distance correlations method of (Székely et al., 2007). (d) Difference between correlation strengths derived in 3(b) and those found at the same locations from 2079-2100. (e) and (f) As per 3(c) and 3(d), but for the equivalently-sized, randomly-sampled training set.

important to note that we should expect these differences to be exaggerated in the real world, where the system is significantly more complex.

Additionally, our results also demonstrate how spatial sampling bias can significantly alter the patterns of apparent relationships between environmental predictors and plankton biomass. The association strengths identified in the measurements-derived sample vary considerably from those found in the random sample of equivalent size (see Figure 3(c) vs. 3(e)). Importantly, this finding is robust across a range of sample sizes, where almost identical patterns of correlations are seen in the 3586 cell case as in the 25,683 cell case, as well as across several methods of deriving correlations (see Figure S12). Nonetheless, the spatial patterns of over and under-prediction derived from the GAMs SDM are not merely the result of spatiotemporal measurement biases. We see remarkable agreement in these broad qualitative patterns between the predictions generated from measurements-derived and random samples ((c) and (d) of Figures 1, and S1–6, and Figures S8 and S9). Ocean measurement biases may explain some element of the tendency towards over-estimation of historical biogeography/abundances; perhaps because measurements have more often been made in places with higher than average abundances. In all cases, training the statistical model on a non-biased dataset reduces the severity of over and under-

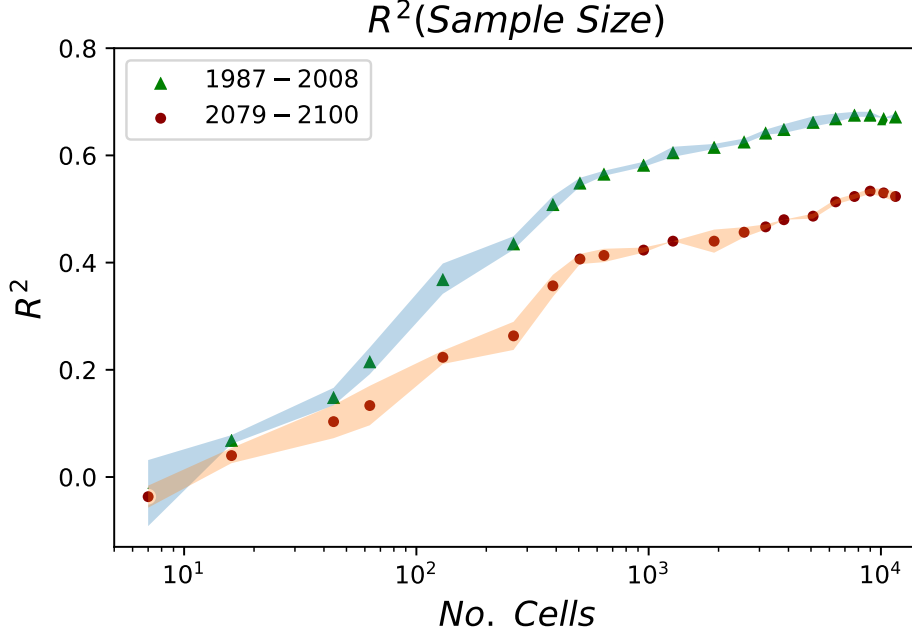


Figure 4: R^2 of GAMs SDM model prediction as a function of sample size. Points are the mean R^2 value for coccolithophore predictions from three independent randomly-generated training sets for each of the 18 sample sizes, ranging from $N=63$ to $N=11,557$. Shading is the standard deviation.

prediction, especially for spatial predictions (Figure S8(e) and S9(e)). But the same broad biogeographic patterns remain, indicating that the SDM is failing to effectively capture changes over time, despite its relatively robust performance according to the broad brush strokes of summary statistics (Figure S4(e) and S4(f)).

The fraction of variance that the SDM can account for saturates with sample size well below 100%, perhaps implying a potential ceiling on predictive ability. Nonetheless, a number of optimisations could be implemented to improve predictive skill; potentially in the SDM developed for the current case study, but certainly in real-world applications. First, we note that an unrepresentative training set presence/absence ratio compared to the population can lead to an unreliable representation of presence/absence in the resulting predictions. To avoid this possibility, researchers working with real observational data will sometimes employ resampling techniques (e.g. (Wei & Dunbrack, 2013)) to account for this effect. By contrast, our experimental design permitted us the unusual opportunity of testing our outcomes alongside a range of representative, randomly-sampled datasets spanning the surface ocean. These unbiased samples are representative of the presence/absence ratios of the population, and thus act as a control for our observations-derived test case. Given the broadly similar patterns of over and underprediction found across test cases, we do not employ resampling techniques here, but we encourage their application in real-world settings.

Related also to the more flexible nature of our study in comparison to correlative SDMs built from real-world observations, is the manner in which we approach training, validation and testing datasets. In some cases, machine learning practitioners working with real-world data, and their associated limitations, might reserve a proportion of the

training set for model validation, as well as an independent, but similarly-distributed, dataset for performance testing. A validation set allows for optimisation via the fine-tuning of model parameters, and for the avoidance of over-fitting, while the test set permits evaluation of model skill. Here, we use whole-ocean Darwin Model output as our test set for evaluating overall performance. Given model response to sensitivity tests, and GAM's natural robustness to over-fitting as a result of predictor function regularisation, we do not explicitly employ a validation set. Model skill could be improved with parameter fine-tuning, especially in the spatial predictions test case. But it is less clear whether fine-tuning for performance using a training set sampled from the Darwin Model ocean of 1987-2008 would improve end-of-century predictions, for reasons that we will return to as this discussion progresses. Additionally, we speculate that our decision to train the GAMs SDM using the entire measurements-derived sample might itself yield improvements relative to splitting the samples into training, testing and validation subsamples.

The median overestimations of the GAMs SDM compared to the Darwin 'ground truth', even when using randomly sampled training data, also implies that these predicted abundance distributions are less skewed than the Darwin model distributions, which are, in turn, less skewed than distributions in the the real ocean. That is not to say, however, that all correlative SDMs will yield equivalent outcomes, regardless of the empirical models at their cores. Recent work by (Rudy et al., 2017) demonstrates that empirical methods can reliably extract the underlying mechanistic equations that govern a dynamical system. Similarly, (Holder & Gnanadesikan, 2021) evaluate random forest (RF) and neural network ensembles (NNE) in their ability to resolve the underlying intrinsic relationships between plankton biomass and environmental predictors, from the apparent relationships in the data. They demonstrate variability in predictive skill across different empirical test cases, and find that NNE's yield overall superior performance; particularly in the case where plankton growth rates respond rapidly to environmental change, as might be expected in many real-world ocean environments. These hybrid methods represent a potential step toward building more skillful and descriptive models.

Although improvements to overall predictive skill might be made through model optimisation techniques, we argue here that the assumptions and uncertainties inherent to correlative SDMs apply fundamental limits to their utility. For instance, although we might feasibly achieve a better fit to the training data, questions still remain as to whether the environmental data included in the model reflect the true and complete niche requirements of the target species'. Even if we were to overcome this issue, using environmental correlates of distribution to predict abundance elsewhere in space and time implies that the distributions in the training data are at equilibrium, such that the niche is 'fully occupied'. This may not be the case, as an otherwise suitable niche for a given species might have experienced some recent perturbation that temporarily reduces its equilibrium population density.

Empirical methods that extract the intrinsic drivers of plankton abundance and distribution (as derived in laboratory settings) might also yield considerable improvements to predictive capabilities of correlative SDMs. If factors such as spatiotemporal sampling bias and spatial autocorrelation in ocean measurements can also be accounted for, predictive skill might be greatly improved, especially in spatial extrapolations. However, appreciable improvements to multidecadal predictions of how plankton communities might respond to climate change would still not be guaranteed; we cannot assume that a specie's niche envelope is fixed and immutable over time. This is clearly demonstrated in our results; but we should expect the predictive skill of correlative SDMs applied to real world data to yield poorer results still. For instance, there are many more degrees of freedom in real-world interactions between plankton individuals, communities, and the wider ecosystem and environment. In addition to the controlling influence of e.g. nutrient supply rate, physical transport processes and level of top down pressure, plankton are also able to adapt genetically, epigenetically and plastically to change. With

their short generation times and high biodiversity, we might expect that even intrinsic relationships could change over the course of a century. This is especially likely in such a dynamic, randomly-perturbed, and far-from-equilibrium environment, where conditions are ideal for unpredictable emergent phenomena to arise. By contrast, all such elements within the Darwin Model are simplified by design, and intrinsic relationships are held steady over time, such that the spatiotemporal variability in apparent relationships seen here are the product of many fewer sources of complexity, right down to how climate change proceeds (a known quantity in the Darwin Model, and yet another significant source of uncertainty in the real world).

We focus here on deriving our SDM using a statistical learning model that, for reasons outlined in Materials & Methods, we believe makes for an excellent case study. Our investigation has allowed us to better clarify the strengths and limitations of such an approach, as applied in the current context. Owing to the complexity and ever-changing nature of the system, some of these limitations could be fundamental and unavoidable, particularly when extrapolating far beyond the training regime.

Methodologically, the broader approach we have presented of applying an empirical model to output from a numerical model may be useful for addressing a number of additional questions. These might include evaluating how best to empirically model whole-ecosystem properties, such as diversity, from observations, or assessing where and when to make new observations to maximise information content about global plankton biogeography. But, as our results here have demonstrated and reinforced, it is important to be aware of the strengths and limitations of this approach, especially when dealing with a high degree of complexity over time.

5 Conclusion

In summary, our results suggest that correlative SDMs like the one developed here can be powerful tools for extrapolating from sparse measurement sets to capture the qualitative spatial patterns of plankton biomass in the present-day ocean. However, their predictions are especially sensitive to the spatiotemporal bias in historical measurements, and can tend towards overprediction if not properly accounted for. In addition, such models demonstrably struggle to predict future plankton biomass because the spatial and temporal complexity of the physical, chemical and biological interactions that characterise the system give rise to a variability that cannot be accurately predicted decades ahead of time from correlations in contemporary data. The changes in relationship between environmental variables and the plankton abundances demonstrated in the current work could be greatly exaggerated in correlative SDMs that tackle the significantly more complex task of predicting real-world plankton biogeography using sparse observational data.

Acknowledgments

Ward acknowledges support from a Royal Society University Research Fellowship. Dutkiewicz acknowledges support from the Simons Collaboration on Computational Biogeochemical Modelling of Marine Ecosystems (CBIOMES)(Grant Id: 549931) and from the United States Air Force Research Laboratory and the United States Air Force Artificial Intelligence Accelerator and was accomplished under Cooperative Agreement Number FA8750-19-2-1000. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the United States Air Force or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes notwithstanding any copyright notation herein. Cael acknowledges support from the National Environmental Research Council (NE/R015953/1) and the Horizon 2020 Framework Programme (820989). The work reflects only the authors' view; the European Commission and their executive agency are not responsible for any use that may be made of the in-

formation the work contains. Finally, the authors' would like to thank the two anonymous reviewers for their insightful comments, which have yielded substantial improvements to the final version of this manuscript.

Code Availability. The physical model used here is available through <http://www.mitgcm.org>, and the generic ecosystem code is available through <http://www.gitlab.com/jahn/gud>. The specific modifications for the setup used here are available via Harvard Dataverse at <http://www.dataverse.harvard.edu/dataverse/>. Note that a more up-to-date version of the ecosystem model used here is available at <http://www.github.com/darwinproject/darwin/>. The code used to process and analyse the data, and to produce the results for this manuscript, is available at <https://github.com/teatauri/stats-biogeo-2021>.

Data Availability. The Darwin Model output used in the current study is available at <http://www.dataverse.harvard.edu/dataverse/>. The dataset will have a doi, and will be hosted through the Harvard Dataverse Darwin project site. The extracted and processed Darwin surface data will also be made similarly available.

References

- Agusti, S., Lubián, L. M., Moreno-Ostos, E., Estrada, M., & Duarte, C. M. (2019). Projected changes in photosynthetic picoplankton in a warmer subtropical ocean. *Front. Mar. Sci.*, 5.
- Barton, A. D., Pershing, A. J., Litchman, E., Record, N. R., Edwards, K. F., Finkel, Z. V., ... Ward, B. A. (2013). The biogeography of marine plankton traits. *Ecology Letters*, 16(4), 522–534.
- Benedetti, F., Vogt, M., Elizondo, U., Righetti, D., Zimmermann, N. E., & Gruber, N. (2021). Major restructuring of marine plankton assemblages under global warming. *Nature Communications*, 12, 5226.
- Cabré, A., Marinov, I., & Leung, S. (2015). Consistent global responses of marine ecosystems to future climate change across the IPCC AR5 earth system models. *Clim. Dyn.*, 45(5), 1253–1280.
- Carvalho, D. V., Pereira, E. M., & Cardosa, J. S. (2019). Machine learning interpretability: A survey on methods and metrics. *Electronics*, 8, 832.
- Cermeño, P., Dutkiewicz, S., Harris, R. P., Follows, M., Schofield, O., & Falkowski, P. G. (2008). The role of nutricline depth in regulating the ocean carbon cycle. *PNAS*, 105(51), 20344–20349.
- Dutkiewicz, S., Boyd, P. W., & Riebesell, U. (2021). Exploring biogeochemical and ecological redundancy in phytoplankton communities in the global ocean. *Global Change Biology*, 27(6), 1196–1213.
- Dutkiewicz, S., Cermeño, P., Jahn, O., Follows, M. J., Hickman, A. E., Taniguchi, D. A. A., & Ward, B. A. (2020). Dimensions of marine phytoplankton diversity. *Biogeosciences*, 17(3), 609–634.
- Dutkiewicz, S., Follows, M. J., & Bragg, J. G. (2009). Modeling the coupling of ocean ecology and biogeochemistry. *Global Biogeochemical Cycles*, 23(4).
- Dutkiewicz, S., Scott, J. R., & Follows, M. J. (2013). Winners and losers: Ecological and biogeochemical changes in a warming ocean. *Global Biogeochemical Cycles*, 27(2), 463–477.
- Dutkiewicz, S., Ward, B. A., Scott, J. R., & Follows, M. J. (2014). Understanding predicted shifts in diazotroph biogeography using resource competition theory. *Biogeosciences*, 11(19), 5445–5461.
- Falkowski, P. G., Barber, R. T., & Smetacek, V. (1998). Biogeochemical controls and feedbacks on ocean primary production. *Science*, 281(5374), 200–206.
- Falkowski, P. G., Fenchel, T., & Delong, E. F. (2008). The microbial engines that drive earth's biogeochemical cycles. *Science*, 320(5879), 1034–1039.

- Fenchel, T. (1988). Marine plankton food chains. *Ann. Rev. Eco. Sys.*, 19(1), 19–38.
- Flombaum, P., Wang, W.-L., Primeau, F. W., & Martiny, A. C. (2020). Global picophytoplankton niche partitioning predicts overall positive response to ocean warming. *Nature Geoscience*, 13(2), 116–120.
- Fuhrman, J. A. (2009). Microbial community structure and its functional implications. *Nature*, 459(7244), 193–199.
- Graff, J., Westberry, T., Milligan, A., Brown, M., Dall’Olmo, G., Reifel, K., & Behrenfeld, M. (2016). Photoacclimation of natural phytoplankton communities. *Marine Ecology Progress Series*, 542.
- Guidi, L., Chaffron, S., Bittner, L., Eveillard, D., Larhlimi, A., Roux, S., ... Gorsky, G. (2016). Plankton networks driving carbon export in the oligotrophic ocean. *Nature*, 532(7600), 465–470.
- Guidi, L., Stemann, L., Jackson, G. A., Ibanez, F., Claustre, H., Legendre, L., ... Gorsky, G. (2009). Effects of phytoplankton community on production, size, and export of large aggregates: A world-ocean analysis. *Limnology and Oceanography*, 54(6), 1951–1963.
- Hastie, T., & Tibshirani, R. (1986). Generalized additive models. *Statistical Science*, 1(3), 297–310.
- Holder, C., & Gnanadesikan, A. (2021). Can machine learning extract the mechanisms controlling phytoplankton growth from large-scale observations? – a proof-of-concept study. *Biogeosciences*, 18, 1941–1970.
- Hutchinson, G. E. (1961). The paradox of the plankton. *Amer. Naturalist*.
- Ibarbalz, F. M., Henry, N., Brandão, M. C., Martini, S., Busseni, G., Byrne, H., ... Zinger, L. (2019). Global trends in marine plankton diversity across kingdoms of life. *Cell*, 179(5), 1084–1097.e21.
- Irwin, A. J., Nelles, A. M., & Finkel, Z. V. (2012). Phytoplankton niches estimated from field data. *Limnology and Oceanography*, 57(3), 787–797.
- Kuhn, A. M., Dutkiewicz, S., Jahn, O., Clayton, S., Rynearson, T. A., Mazloff, M. R., & Barton, A. D. (2019). Temporal and spatial scales of correlation in marine phytoplankton communities. *Journal of Geophysical Research: Oceans*, 124(12), 9417–9438.
- Lombard, F., Boss, E., Waite, A. M., Vogt, M., Uitz, J., Stemann, L., ... Appeltans, W. (2019). Globally consistent quantitative observations of planktonic ecosystems. *Front. Mar. Sci.*, 6.
- Marinov, I., Doney, S. C., & Lima, I. D. (2010). Response of ocean phytoplankton community structure to climate change over the 21st century: partitioning the effects of nutrients, temperature and light. *Biogeosciences*, 7(12), 3941–3959.
- Marinov, I., Gnanadesikan, A., Sarmiento, J. L., Toggweiler, J. R., Follows, M., & Mignone, B. K. (2008). Impact of oceanic circulation on biological carbon storage in the ocean and atmospheric pCO₂. *Global Biogeochemical Cycles*, 22(3).
- Martiny, A., & Flombaum, P. (2020). Global observations prochlorococcus, synechococcus, and picoeukaryotic phytoplankton with ancillary environmental data from 1987 to 2008.
- Melo-Merino, S. M., Reyes-Bonilla, H., & Lira-Noriega, A. (2020). Ecological niche models and species distribution models in marine environments: A literature review and spatial analysis of evidence. *Ecological Modelling*, 415, 108837.
- Prinn, R. G. (2013). Development and application of earth system models. *PNAS*, 110, 3673–3680.
- Righetti, D., Vogt, M., Gruber, N., Psomas, A., & Zimmermann, N. E. (2019). Global pattern of phytoplankton diversity driven by temperature and environmental variability. *Science Advances*, 5(5), eaau6253.
- Rudy, S. H., Brunton, S. L., Proctor, J. L., & Kutz, J. N. (2017). Data-driven discovery of partial differential equations. *Science Advances*, 3, e1602614.

- 626 Rutherford, S., D’Hondt, S., & Prell, W. (1999). Environmental controls on the geo-
627 graphic distribution of zooplankton diversity. *Nature*, 400(6746), 749–753.
- 628 Servén, D., & Brummitt, C. (2018). pygam: Generalized additive models in python.
629 *Zenodo*.
- 630 Sokolov, A. (2005). The MIT integrated global system model (IGSM) version 2:
631 Model description and baseline evaluation. , 46.
- 632 Sonnewald, M., Dutkiewicz, S., Hill, C., & Forget, G. (2020). Elucidating ecological
633 complexity: Unsupervised learning determines global marine eco-provinces.
634 *Science Advances*, 6(22), 2375–2548.
- 635 Székely, G. J., Rizzo, M. L., & Bakirov, N. K. (2007). Measuring and testing depen-
636 dence by correlation of distances. *The Annals of Statistics*, 35(6), 2769–2794.
- 637 Tang, W., & Cassar, N. (2019). Data-driven modeling of the distribution of di-
638 azotrophs in the global ocean. *Geophysical Research Letters*, 46(21), 12258–
639 12269.
- 640 Tittensor, D. P., Mora, C., Jetz, W., Lotze, H. K., Ricard, D., Berghe, E. V., &
641 Worm, B. (2010). Global patterns and predictors of marine biodiversity across
642 taxa. *Nature*, 466(7310), 1098–1101.
- 643 Ward, B. A., Dutkiewicz, S., & Follows, M. J. (2014). Modelling spatial and tem-
644 poral patterns in size-structured marine plankton communities: top-down and
645 bottom-up controls. *Journal of Plankton Research*, 36(1), 31–47.
- 646 Wei, Q., & Dunbrack, R. (2013). The role of balanced training and testing data sets
647 for binary classifiers in bioinformatics. *PloS One*, 8, e67863.
- 648 Wiens, J. A., Stralberg, D., Jongsomjit, D., Howell, C. A., & Snyder, M. (2009).
649 Niches, models, and climate change: Assessing the assumptions and uncertain-
650 ties. *PNAS*, 106, 19729–19736.