

Using Machine Learning to improve the representation of Phytoplankton Dynamics

in Earth System Models

Sandupal Dutta, Anand Gnanadesikan (sdutta12@jh.edu) Department of Earth & Planetary Sciences, Johns Hopkins University



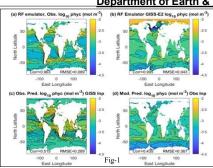
Fig-1: Plankton bloom over Gulf of

Kostadinov: Phys

Introduction.

Phytoplankton generate 50% of the oxygen produced yearly, fix approx 50 Gigatons Carbon per year, Phytoplankton size classes (PSCs) are picoplankton (<2µm), nanoplankton (2-20µm) and microplankton (>20µm). PSCs are closely related to photosynthetic efficiency, sinking rate, and structure of the marine food chain affecting the biological pump.

Hirata: Chlor-a

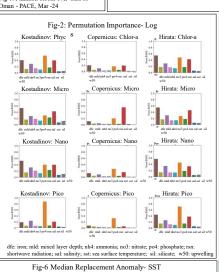


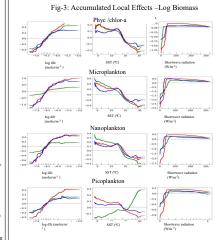
Random Forest (RF). RF predicts PSCs for Kostadinov, MODIS+ Hirata and Copernicus GlobColor very well. Metrics used are Normalized Root Mean Squared Error and R2 Score. Normalized RMSE are in the range of 0.27-0.35, R2 scores are in range between 0.87-0.94.

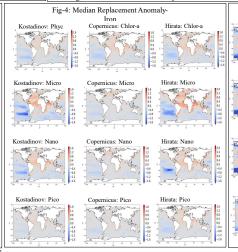
Sources of Errors in ESM. Fig-1 shows distributions of log10 PHYC predicted using RFs. (a) Using observed predictors and the RF emulator trained on observations (Kostadinov et al., 2016), Correlation is very high and RMSE is very small. (b) Using modeled predictors and the RF emulator trained on the model. Correlation is again very high and RMSE very small. (c) Using modeled predictors and the RF trained on the observations. Decline in prediction shows impact of biases in predictors. (d) Using observed predictors and the RF trained on the model. Decline in prediction shows impact of the modeled apparent relationships differing from observed relationships.

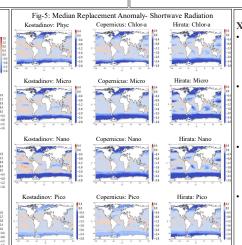
Explainable Machine Learning.

- · Permutation Importance Analysis: Measures importance of a feature as increase in prediction error after permuting the feature's
- Median Replacement Anomaly: Replaces value of one predictor with its median observed value. Difference with original prediction gives the information on spatio-temporal variation of
- Accumulated Local Effects (ALE): Similar to the Partial Dependence Plots (PDP), aiming to describe how features influence average model predictions. However, ALE addresses the bias that arises in the PDP when the feature of interest is correlated with other features.









XAI Results.

- Shortwave radiation (rsn), Sea Surface Temperature (sst), (dfe), Ammonia (nh4) important (Fig.-5).
- Shortwave radiation more important at higher latitudes and exerts influence in a narrow range (30-40 W/m-2) (Fig-5, right-hand column of Fig.- 3).
- SST has negative influence mostly over tropics, mid-latitude gyres and positive influence over higher latitudes. (Fig.- 6)
- Iron limitation affects all PSCs till about 2 NM range but sensitivity varies across PSCs with microplankton showing highest and picoplantkon the lowest. (left-hand column of Fig. 3)

Plankton Dynamics

$$\frac{\partial P_i}{\partial t} = \mu_{ref} * T_f * L_i * V_i * P_i - g_z^{max} * Z * \frac{P_i}{K_{P_i} + P_i} - m_i * P_i - \alpha_i * P_i^{1.75}$$

$$L_i = 1 - exp\left(\frac{-\alpha_i^{Chl} * \theta_i^{C} * I}{\mu_{ref} * T_f * V_i}\right) \qquad \qquad V_i = min\left(\frac{N_j}{K_{ji} + N_j}\right)$$

$$1.7 \frac{T - 30^{\circ}C}{10^{\circ}C}$$
 Fig-

SciML.

(a) Since, observated time series datasets of nutrients of the ocean are sparse and not available with adequate spatial and temporal resolution, hence fully data-driven approaches may not feasible. Therefore, the initial endeavor would be to use gray box modelling to find optimal system parameters for the plankton dynamics.

Gray box models blend the advantages of both white and black box approaches. They are capable of incorporating known physics of plankton dynamics (white box) into machine learning techniques (black box) thereby using data-driven methods to tune parameters of a physically based model. The equations of Plankton dynamics encoded in MARBL is shown in Fig.-7 (Long et al., 2021).

Initial approach would be to use Universal Differential Equation (UDE) formulation with the MARBL equations to construct datadriven gray box representations of the Plankton dynamics. The endeavor is to learn parameters such as μ_{ref} , m_i , α_i using observations. The idea is to replace unknown parts of the Plankton dynamics equation from MARBL with neural networks and learn the optimal parameters using data-driven approaches.

(b) The next step would be to undertake "equation discovery", or "data-driven discovery of partial differential equations" or "learning hidden physics". This involves identifying the underlying mathematical models (governing equations) from observed data, especially when knowledge about the system is incomplete. Frameworks such as Deep Hidden Physics Operator (DHPO) - Discovering Physics using DeepONet will be explored at this stage (Kag et al., 2024). Conceptually, it aims to discover an unknown physics operator N in a general nonlinear PDE of the form

$$\frac{\partial P_i}{\partial t} = \mathbb{N}(\sum P_i, Dr, Nut)$$

A custom loss function can be used to train such an architecture. A possible formulation would be to combine the data loss and the PDE loss term into a combined Sum of Squared Errors Loss (SSE) as shown below (Raissi, 2018).

$$SSE = \sum (|P_i - \hat{P}_i|^2 + |\frac{\partial P_i}{\partial t} - \mathbb{N}(\sum P_i, Dr, Nut)|^2)$$

(c) Subsequently, the approach would be to use symbolic regression (e.g., Sparse Identification of Nonlinear Dynamics, SINDy), LAGRAMGE or other equation discovery systems to extract interpretable forms like

$$\frac{\partial P_i}{\partial t} = \mu(\sum P_i, Dr, Nut) * P_i - m * P_i$$

Note: Work supported by NOAA grant NA21OAR4310256 and DOE grant SC0025209.

