# Derivation of Marine Inherent Optical Properties
## *A Bayesian Approach*

## Erdem M. Karaköylü & Susanne E. Craig
NASA Ocean Biology Processing Group

erdem.m.karakoylu@nasa.gov — (301) 286-0501

## Introduction

The advent of satellite oceanography has generated tremendous insight into marine biogeophysical processes by providing a global view of phytoplankton distribution dynamics, with measurable impact on various fields like ecology, fisheries, climate science. One of the principle obstacles in this endeavor is that the water's contribution to the sensed light field is dwarfed by that of the atmosphere. Some of this, like Rayleigh scattering, is straightforward to correct for. However, light redirected to the sensor by spatially and temporally variable aerosol distribution has so far been addressed by computing a modelled approximation. This approximation works well in the open ocean, but runs into trouble in coastal regions, where both marine and atmospheric layers are often optically complex, and as yet not well understood..

Here, we propose circumventing the atmospheric complexity of coastal areas by using top-of-the-atmosphere radiance (**TOA**) along with some additional ancillary input. We develop and compare alternative models to estimate phytoplankton absorption (**aph**) - a proxy for phytoplankton distribution - using a bayesian modeling framework.
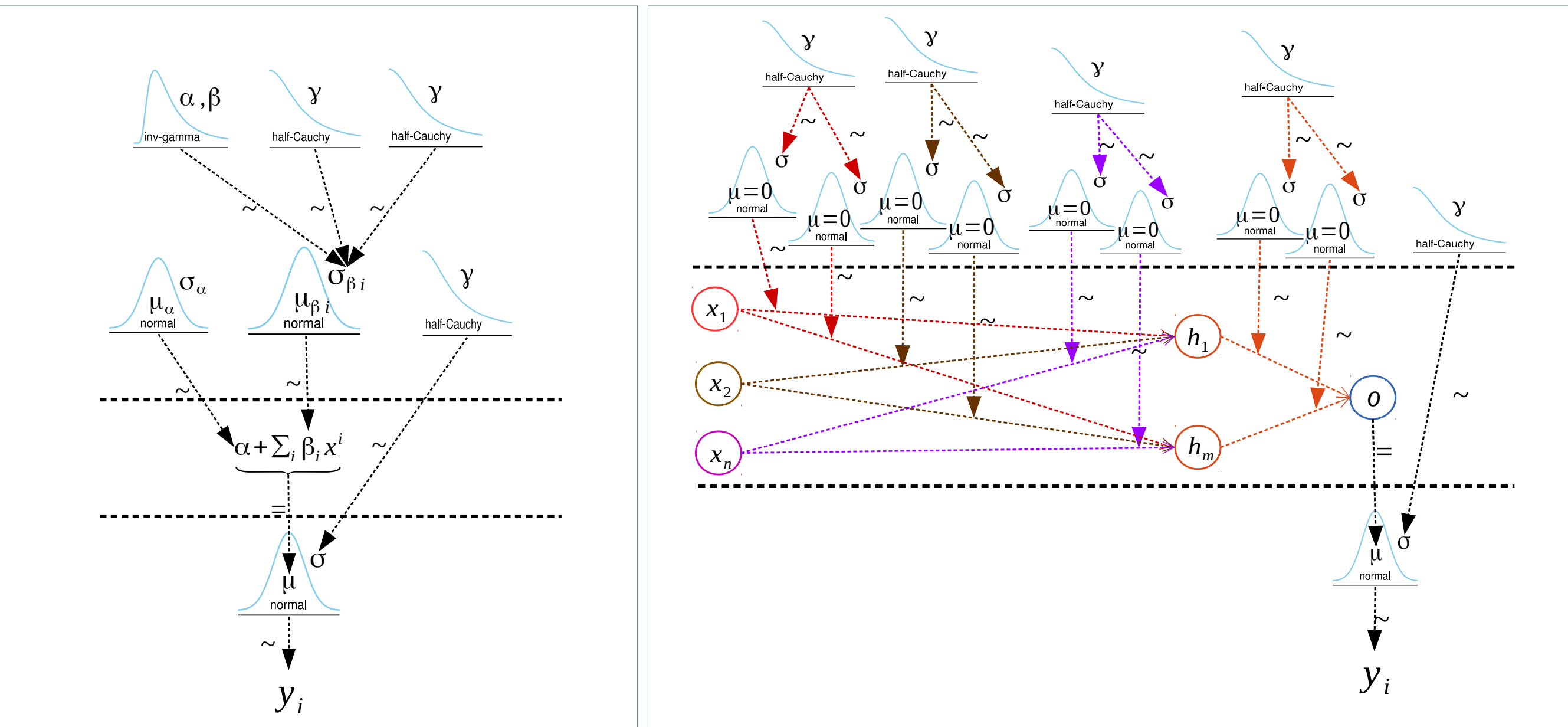
## Objectives

1. Estimate phytoplankton distribution in coastal regions by modeling **aph**.

2. Develop feature selecting *Bayesian linear and nonlinear models*.

3. Use performance assessment and information theory for evaluation and model selection.

## Materials and Methods

### Model Development

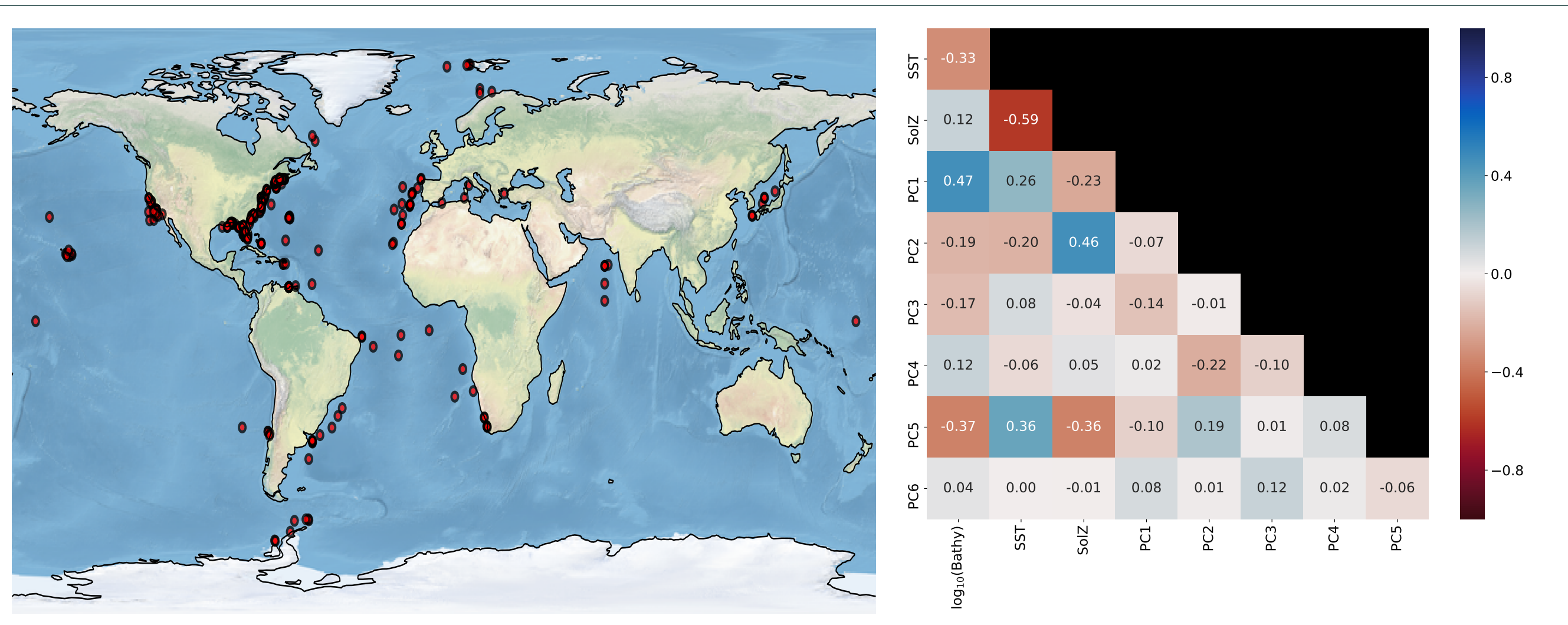Three models were developed using the python probabilistic framework PyMC3[1].

1. Linear regression with regularized horseshoe prior (**rHSP**)[2]; hereafter, **Model 1**.

2. Linear regression with **rHSP** and $1^{st}$ order feature interactions; hereafter, **Model 2**.

3. Bayesian neural network [3]; hereafter, **Model 3**.



**Figure 1: Model Structure.** Horizontal lines separate three conceptual groups; top → priors, middle → likelihood, bottom → outcome distribution. **Left:** Regression with horseshoe priors (Models 1 & 2). **Right:** Bayesian neural network (Model 3). Models shown here are hierarchical, built for automatic feature relevance determination.
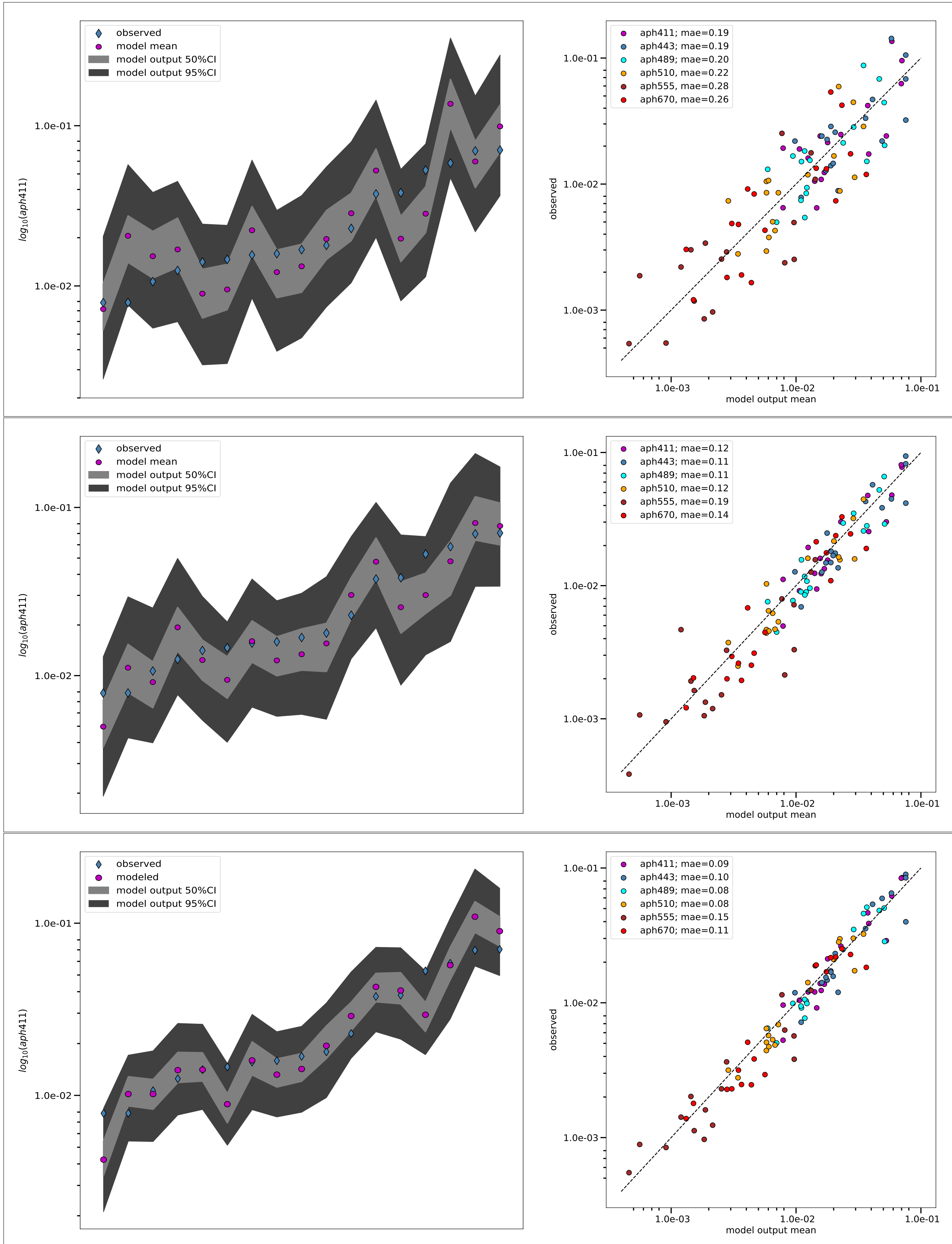
### Data Pre-Processing

Data consisted in satellite **TOA** observations matched up, according to [4], *in-situ* **aph** measurements at 6 wavelengths; 411, 443, 489, 510, 555, and 670 nm. Pre-processing highlights: **(1)** principle components (**PC**) computed from correlated TOA radiance; **(2)** features include PCs, water temperature (**SST**), solar zenith angle (**SolZ**), depth (**Bathy**); **(3)** although not strictly required in a Bayesian setting, data was split into training/testing (out-of-sample) sets.



**Figure 2: Data Overview. Left**: In-situ sampling locations are mostly coastal; continental or insular. **Right**: Pairwise Pearson correlation of features used in all models.

## Results



**Figure 3: Model Evaluation on Test Set.** Top, middle and bottom panels correspond to Models 1, 2, and 3, respectively. **Left**: Out-of-sample observations of **aph** at 411nm, in relations to model posterior predictive mean, 50%, and 95% credibility interval. **Right**: Out-of-sample observations against model predictions, for all bands of **aph**, with goodness-of-fit scored by mean absolute error, (**mae**).

| | WAIC | pWAIC | dWAIC | weight | SE | dSE |
|---|---|---|---|---|---|---|
| **Model 3** | **-183.21** | **29** | **0** | **0.98** | **21.51** | **0** |
| Model 2 | -70.48 | 32.7 | 112.72 | 0 | 18.49 | 16.09 |
| Model 1 | -25.91 | 9.81 | 157.3 | 0.02 | 15.01 | 19.93 |

**Table 1: W**idely **A**vailable **I**nformation **C**riterion for models predicting **aph** at 411 nm. WAIC takes into account model flexibility and the posterior distribution of a model to predict its performance on future (out-of-sample) data. **WAIC**: lower score predicts better performance; **pWAIC**: effective number of parameters - a measure of model flexibility ; **dWAIC**: difference with lowest WAIC ;**weight**: can be used when ensemble averaging similarly scored models when no clear winner is available; **SE**: standard error of WAIC estimate ; **dSE**: standard error of dWAIC. Here, **Model 3**, the Bayesian Neural Network model is predicted to be a more robust model, and proposed as sole model to be selected.

## Conclusions

- Bayesian inference provides a principled modeling framework.
- Assumptions are explicit resulting in criticizable models that can be built to be comparable.
- Models 1 and 2 selected PCs 1-3, SST, SolZ, and Bathy. as significant features (results not show).
- Model 3 selected PCs 1-3 as significant features (results not show).
- By all measures, **Model 3** is predicted to be the better performing alternative.

## Problems and Opportunities

- Symmetry issues complicates convergence in neural networks.
- More *in-situ* data collection needed for more robust model building.
- Integration of this and other approaches into existing production systems.

## References

[1] John Salvatier, Thomas V. Wiecki, and Christopher Fonnesbeck. Probabilistic programming in python using PyMC3. *PeerJ Computer Science*, 2:e55, 2016.

[2] J. Piironen and A. Vehtari. Sparsity information and regularization in the horseshoe and other shrinkage priors. *Electronic Journal of Statistics*, 11:5018–5051, 2017.

[3] R. M. Neal. *Bayesian Learning for Neural Networks*. Springer, 1996.

[4] S. W. Bailey and P. J. Werdell. A multi-sensor approach for the on-orbit validation of ocean color satellite data products. *Remote Sensing of Environment*, 102(1-2):12–23, 2006.