

## **Machine Learning Estimates of Global Marine Nitrogen Fixation**

**Weiyi Tang<sup>1</sup>, Zuchuan Li<sup>1</sup>, Nicolas Cassar<sup>1,2,\*</sup>**

<sup>1</sup>Division of Earth and Ocean Sciences, Nicholas School of the Environment, Duke University, Durham, NC 27708, USA

<sup>2</sup>Laboratoire des Sciences de l'Environnement Marin (LEMAR), UMR 6539 UBO/CNRS/IRD/IFREMER, Institut Universitaire Européen de la Mer (IUEM), Brest, France

\*Correspondence to: Nicolas Cassar ([Nicolas.Cassar@duke.edu](mailto:Nicolas.Cassar@duke.edu))

### **Contents of this file**

Text S1  
Figures S1 to S7

### **Additional Supporting Information (Files uploaded separately)**

Captions for Movies S1 to S2

24 **Text S1.**

## 25 **Caveats, uncertainties and future improvements**

### 26 **1. Sparse and uneven distribution of observations, and mismatch of N<sub>2</sub> fixation** 27 **observations with predictors**

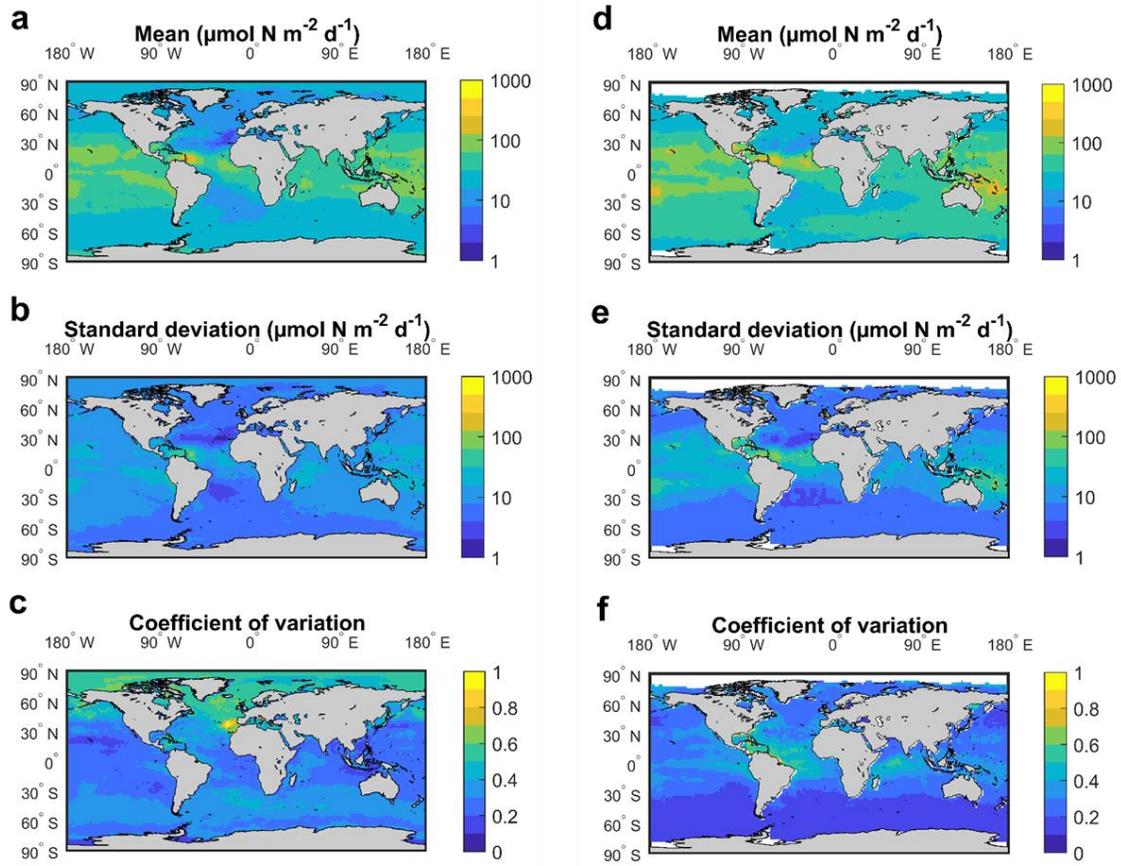
28 While our training dataset includes new observations in the South Pacific, the Indian and  
29 Arctic Oceans (which were not available in Luo et al. (2014)), the majority of  
30 observations remain in the North Atlantic Ocean (Figure 1). For example, only 6 points  
31 are available in the Indian Ocean. This uneven and sparse distribution of observations  
32 may bias the statistical models by giving some regions more weight. Given that factors  
33 regulating N<sub>2</sub> fixation likely vary between biomes and regions (Monteiro, Dutkiewicz, &  
34 Follows, 2011; Ward, Dutkiewicz, Moore, & Follows, 2013; Weber & Deutsch, 2014),  
35 models based on data mainly collected in the warm and oligotrophic waters of the  
36 Atlantic and Pacific Ocean may not accurately represent the other regions, including the  
37 recently discovered niches of N<sub>2</sub> fixation in cold and/or nutrient-rich waters. In addition,  
38 spatial and temporal mismatch, and the coarse spatial and temporal resolution of our  
39 predictors and predictand may introduce noise. Some of the predictors may also work  
40 over longer timescales or larger spatial scales than the ones captured by the short-term  
41 incubations. Finally, some environmental factors like nutrient supply ratios (Ward et al.,  
42 2013) may be better predictors of the presence or absence of diazotrophs rather N<sub>2</sub>  
43 fixation rates. Overall, observations over broader swaths of the oceans will help further  
44 refine the biogeography and magnitude of marine N<sub>2</sub> fixation.

45

### 46 **2. Difference in methods measuring N<sub>2</sub> fixation**

47 Our models rely on measurements of N<sub>2</sub> fixation collected by three different incubation  
48 methods (i.e. ARA, <sup>15</sup>N<sub>2</sub> gas addition, and dissolved <sup>15</sup>N<sub>2</sub> addition). The uncertainties,  
49 assumptions and drawbacks of each method introduce biases and noise in our predictions.  
50 In line with this, training the algorithms with each individual method leads to significantly  
51 different biogeographies of N<sub>2</sub> fixation (Figure S2). The ARA method detects bulk N<sub>2</sub>  
52 fixation rates including the particulate and dissolved products (Mulholland, 2007), but  
53 suffers from variable conversion stoichiometry of acetylene reduction to N<sub>2</sub> fixation  
54 (Wilson, Böttjer, Church, & Karl, 2012) and other issues presented in Cassar et al. (2018).  
55 It also displays a geographical bias with most applications being conducted in the North  
56 Atlantic. The <sup>15</sup>N<sub>2</sub> gas addition method is the most commonly used method so far.  
57 Unfortunately, it has been shown to underestimate N<sub>2</sub> fixation rates because of incomplete  
58 gas-liquid equilibration of the <sup>15</sup>N<sub>2</sub> tracer (Mohr, Großkopf, Wallace, & LaRoche, 2010)  
59 and other issues (Bombar, Paerl, Anderson, & Riemann, 2018; Mulholland, 2007). While  
60 a correction may be applied for the incomplete equilibration (Böttjer et al., 2016), it comes  
61 with significant uncertainty because of variability in the degree of disequilibrium between  
62 studies. The dissolved N<sub>2</sub> addition method is now believed to give the best estimates of *in-*  
63 *situ* N<sub>2</sub> fixation rates. However, the measurements are too few at this time to meaningfully  
64 train our machine learning algorithms (Figure 1). Finally, varying depths of integration  
65 may also lead to significant uncertainties. Some studies report N<sub>2</sub> fixation rates integrated  
66 over the euphotic zone while others report rates to a specific depth (e.g. 200 m). The recent  
67 discovery of aphotic N<sub>2</sub> fixation (Fernandez, Farías, & Ulloa, 2011; Hamersley et al., 2011)  
68 exacerbates this issue. Although rates of N<sub>2</sub> fixation are low at depth, they may be

69 significant when integrated over deep water columns. Observations should therefore be  
70 reported to a depth relevant to N<sub>2</sub> fixation to simplify inter-study comparisons.

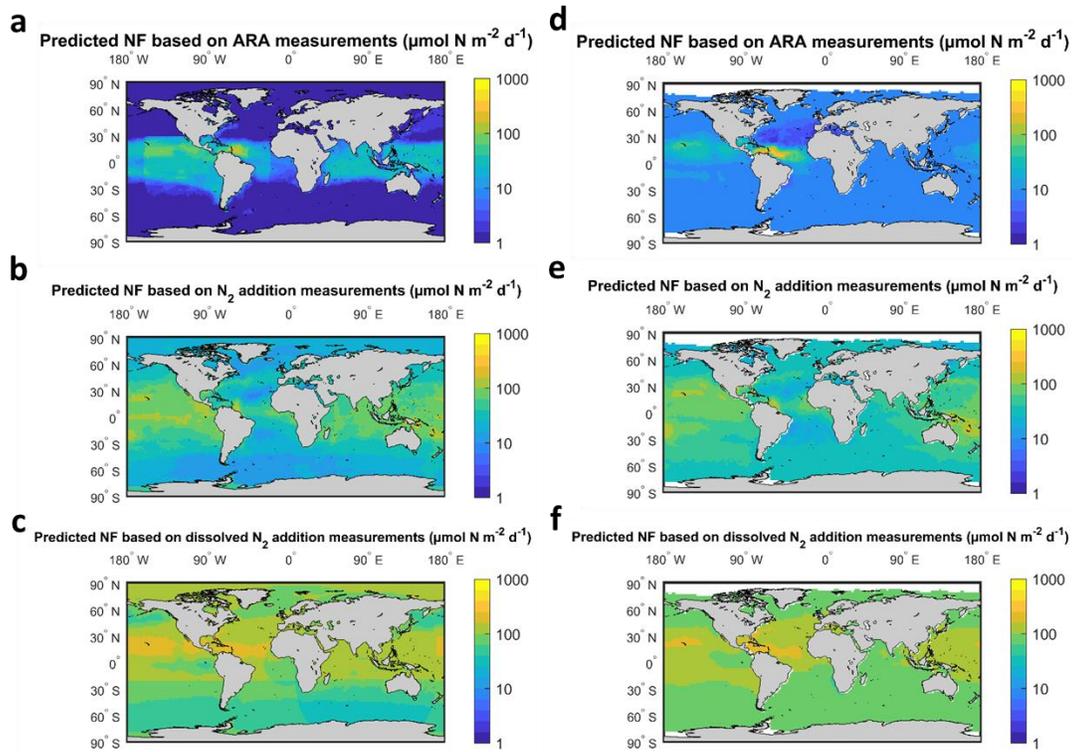


71

72 **Figure S1.** Mean, standard deviation and coefficient of variation of global N<sub>2</sub> fixation from 100  
 73 bootstrap reconstructed N<sub>2</sub> fixation datasets by random forest (RF, a-c), support vector regression  
 74 (SVR, d-f) respectively.

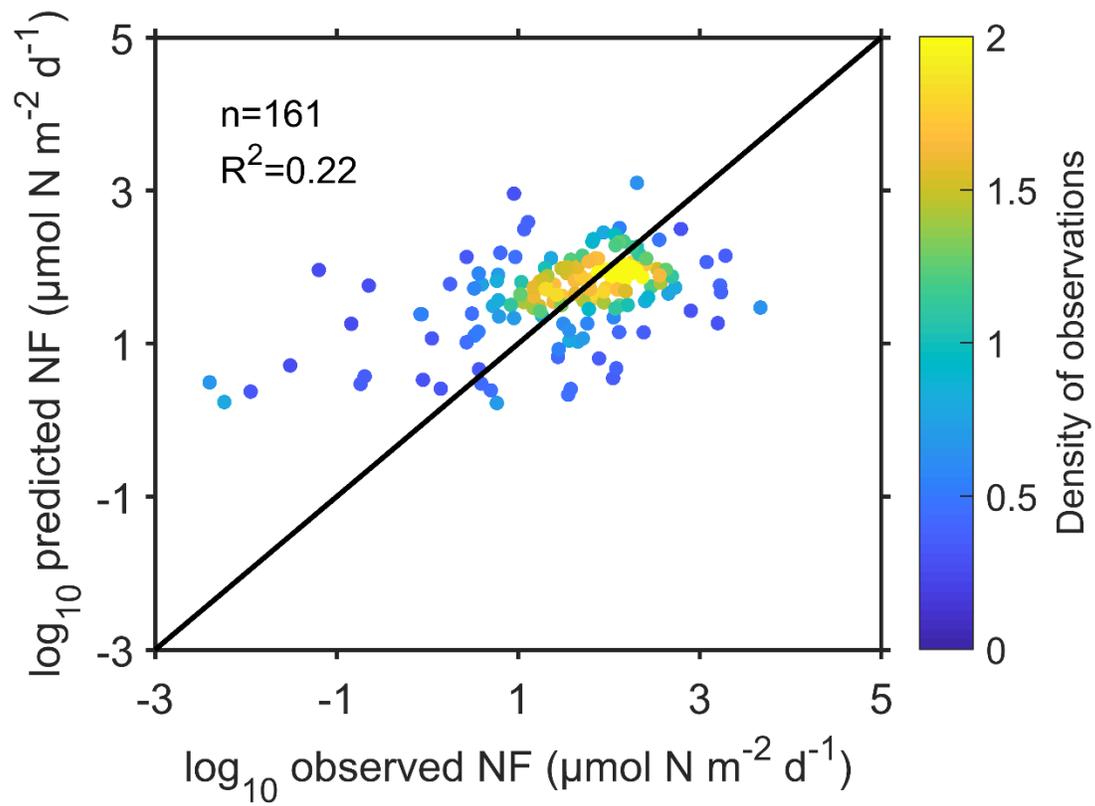
75

76



77

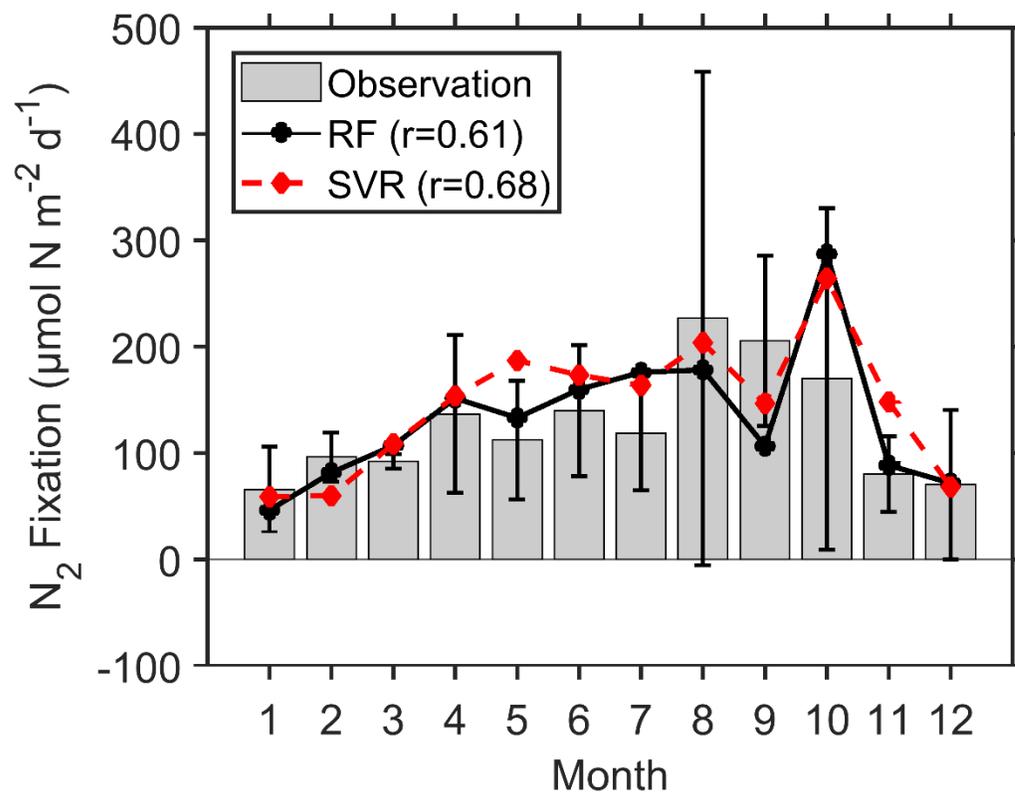
78 **Figure S2.** RF (a-c) and SVR (d-f) model predictions of N<sub>2</sub> fixation based on observations collected  
 79 with a single method. a, d. Acetylene reduction assay (ARA). b, e. <sup>15</sup>N<sub>2</sub> gas addition. c, f. Dissolved  
 80 <sup>15</sup>N<sub>2</sub> addition.



81

82 **Figure S3.** Observed versus simulated  $\text{N}_2$  fixation rates by stepwise multiple linear regression.

83

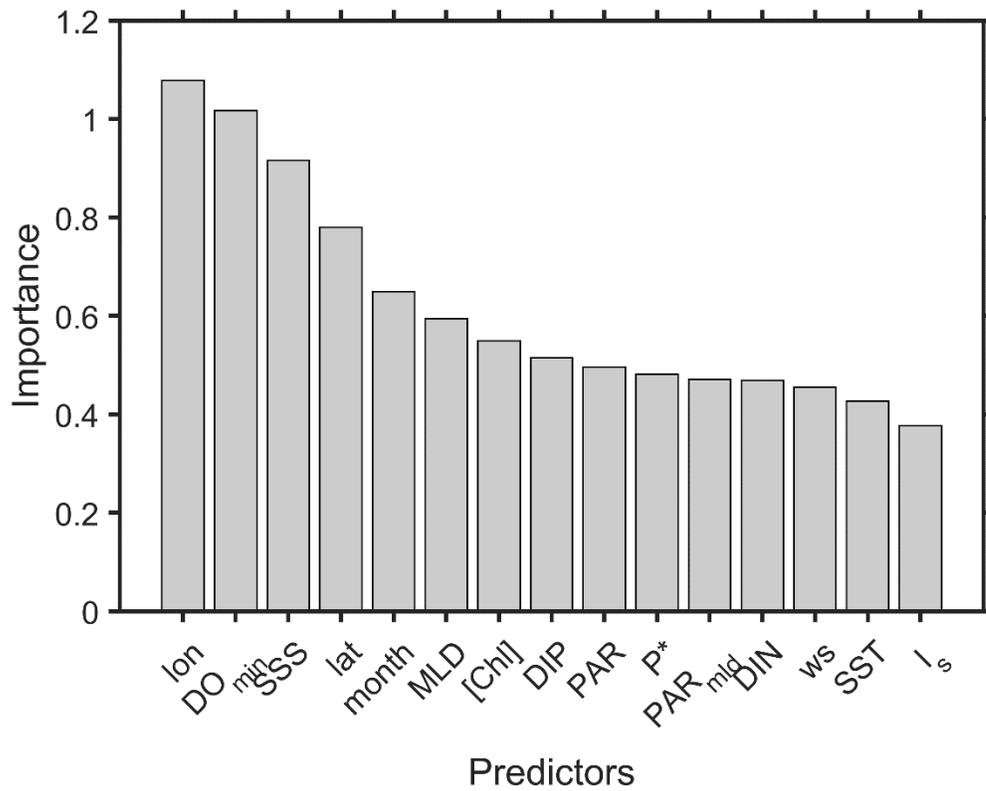


84

85 **Figure S4.** Comparison of observed and modeled seasonal changes in N<sub>2</sub> fixation rates at Hawaii

86 Ocean Time-series (HOT). Bar plot with error bars represents the observed monthly climatology

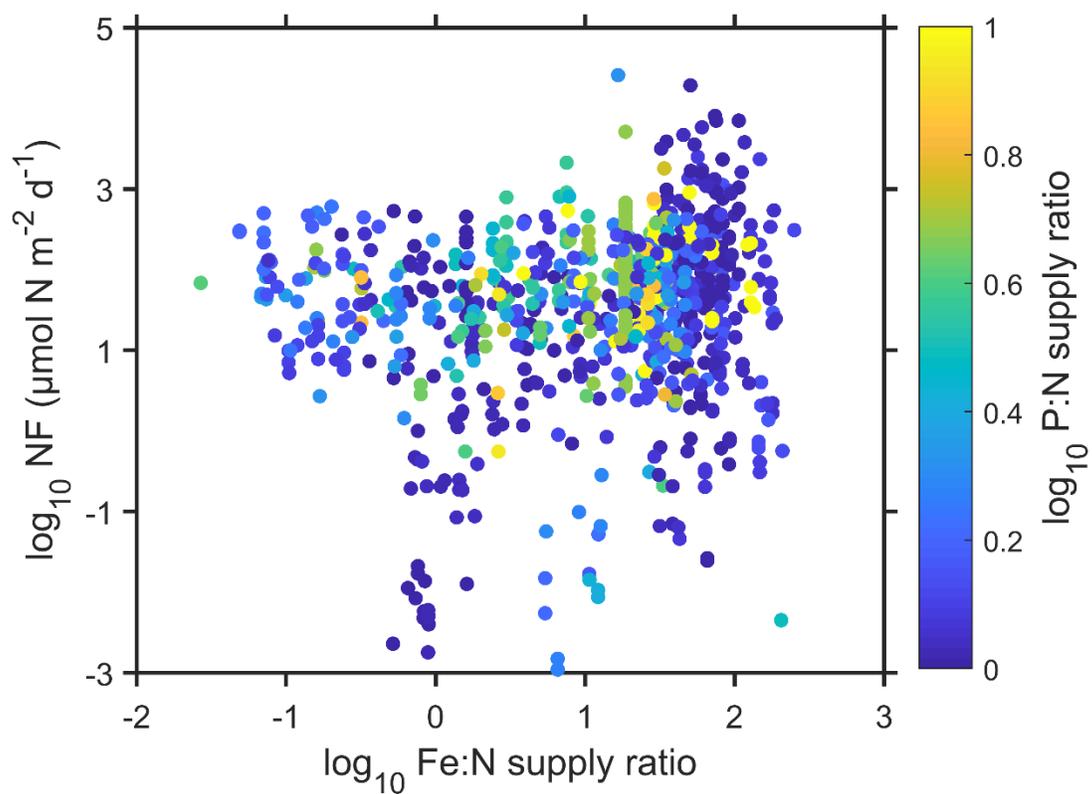
87 of N<sub>2</sub> fixation rates ± one standard deviation at HOT.



88

89

**Figure S5.** Predictor feature importance in random forest.

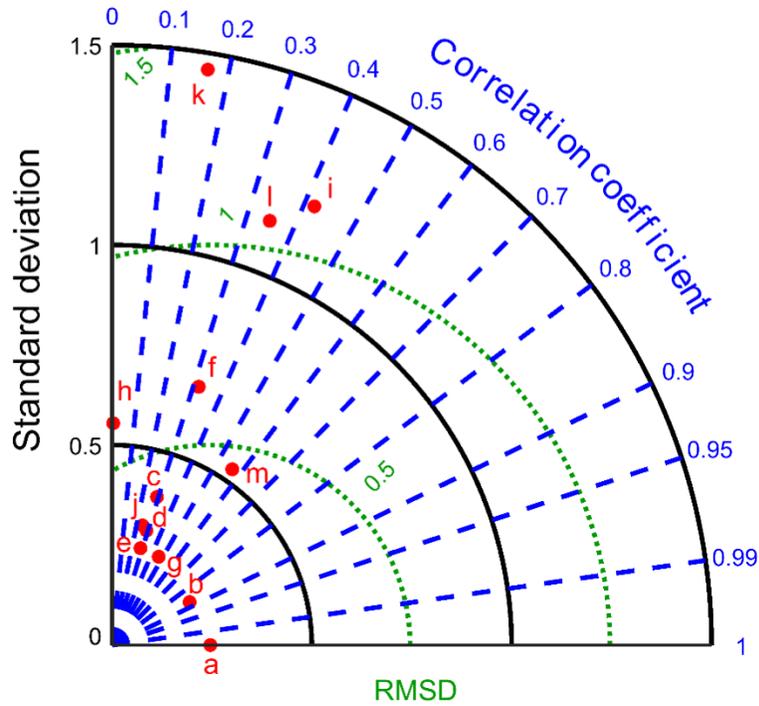


90

91 **Figure S6.** N<sub>2</sub> fixation rates versus Fe:N and P:N supply ratios. Fe:N and P:N supply ratios are

92 from Ward et al. (2013).

93



94

95 **Figure S7.** Taylor diagram of  $N_2$  fixation rates (logarithmic scale) estimated by different models  
 96 with the alphabetical order shown in Figure 4, with the estimate by RF (a) as the reference value.  
 97 Dashed blue and dotted green lines represent the correlation and centered root-mean-square  
 98 difference (RMSD) between estimates by RF and other models, respectively. Solid black lines, the  
 99 radial distance from the origin, represent the standard deviation of the spatial distribution estimated  
 100 by each model, with lower values indicating less spatial variability.

101

102

103

104

105

106 **Movie S1.** Monthly changes of predicted  $N_2$  fixation rates by random forest over the global ocean.

107 **Movie S2.** Monthly changes of predicted  $N_2$  fixation rates by support vector regression over the

108 global ocean.