



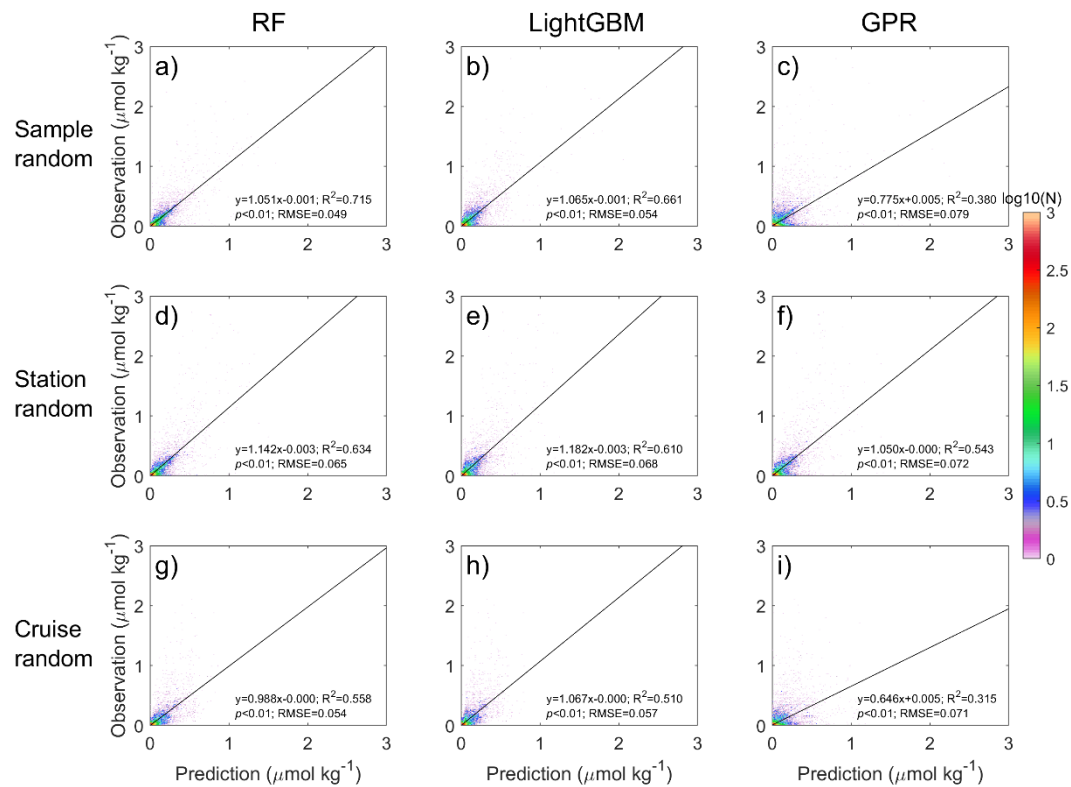
Supplement of

**A historical nutrient dataset (1895–2024) for the North Pacific:
reconstructed from machine learning and hydrographic observations**

Chuanjun Du et al.

Correspondence to: Chuanjun Du (cjdu@hainanu.edu.cn) and Xiaolin Li (xlli@xmu.edu.cn)

The copyright of individual parts of the supplement might differ from the article licence.

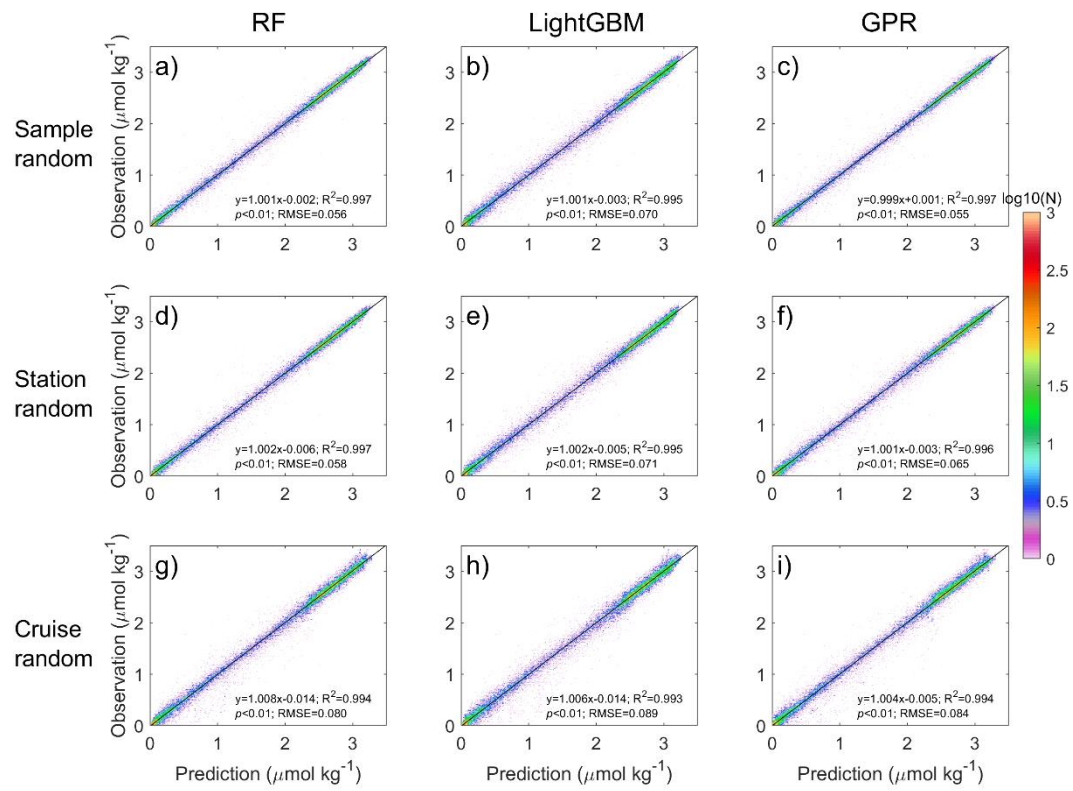


7

8 **Figure S1.** Validating the reconstructed NO_2^- concentrations using leave-one-out cross-
 9 validation with different data selection strategies and machine learning methods. Plots
 10 shown in row 1 correspond to the sample random strategy (a-c), row 2 correspond to
 11 the station random strategy (d-e), and row 3 correspond to the cruise random
 12 strategy (g-i). Plots shown in column 1 correspond to the Random Forest (RF; a, d, and
 13 g), column 2 correspond to the LightGBM (b, e, and h), and column 3 correspond to
 14 the Gaussian Process Regression (GPR; c, f, and i). The black lines and text show the
 15 fitted linear regressions, regression equations, coefficient of determination (R^2), p
 16 values, and Root Mean Squared Errors (RMSE). The color represents the data density
 17 (N , number of observations). Note that the logarithmic scale of N is applied.

18

19



20

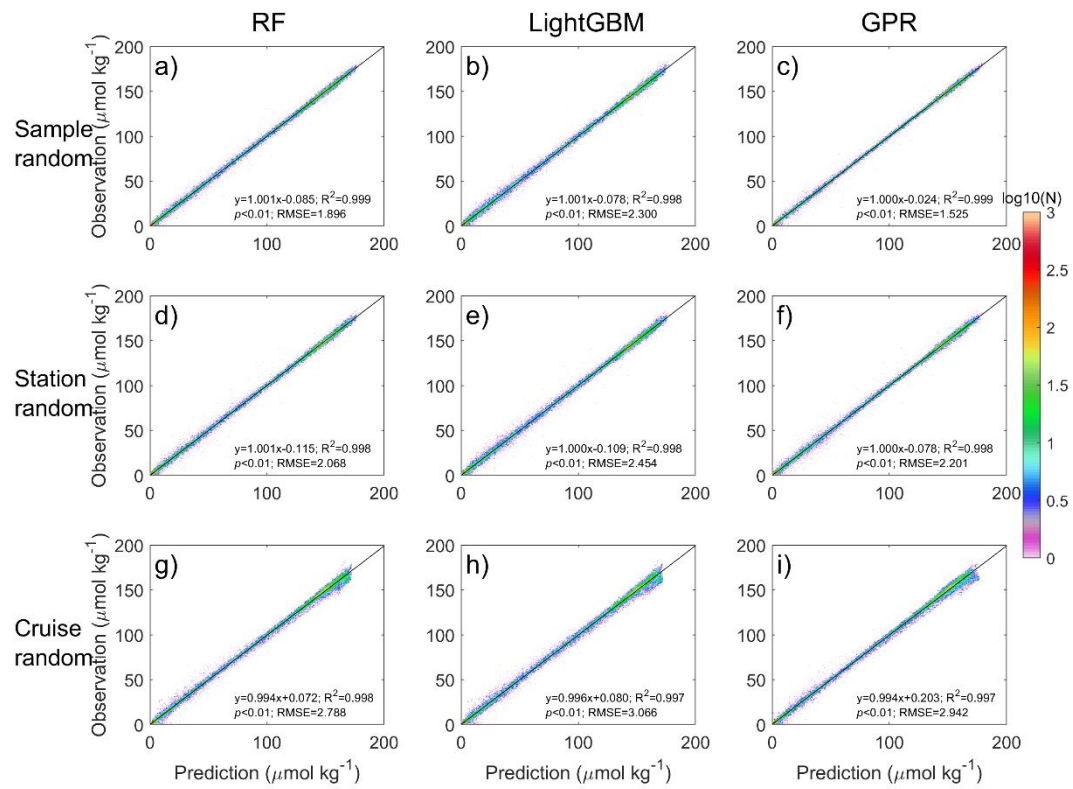
21 **Figure S2.** Similar to Fig. S1, but for DIP.

22

23

24

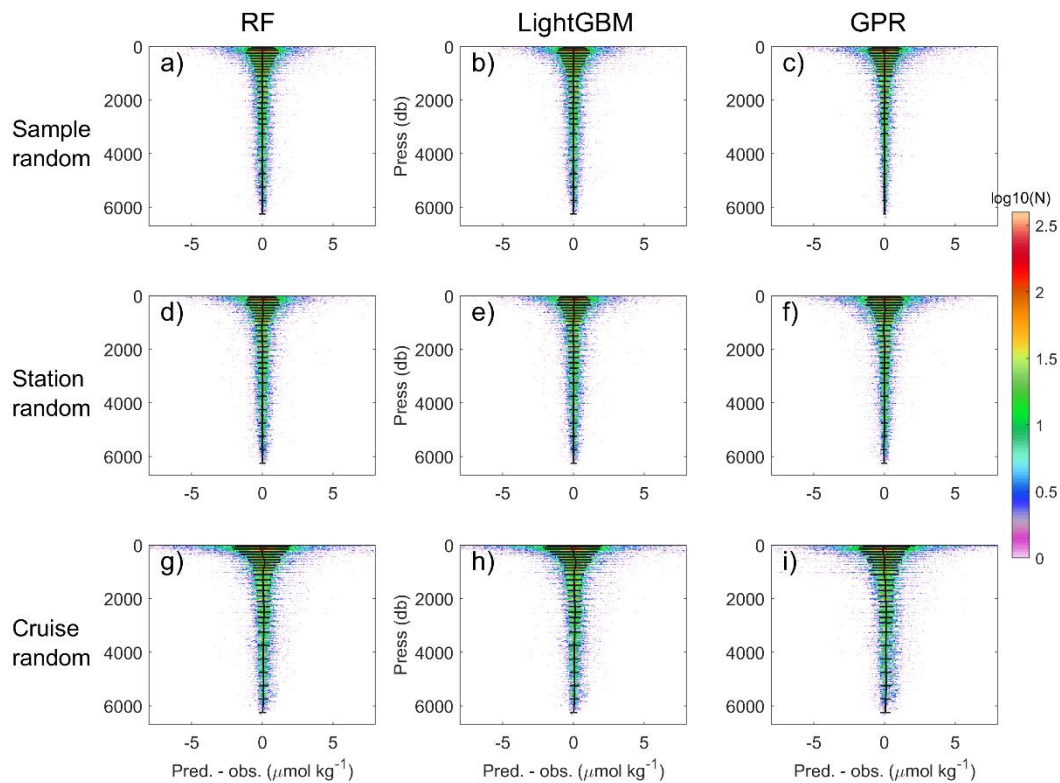
25



26

27 **Figure S3.** Similar to Fig. S1, but for Si(OH)_4 .

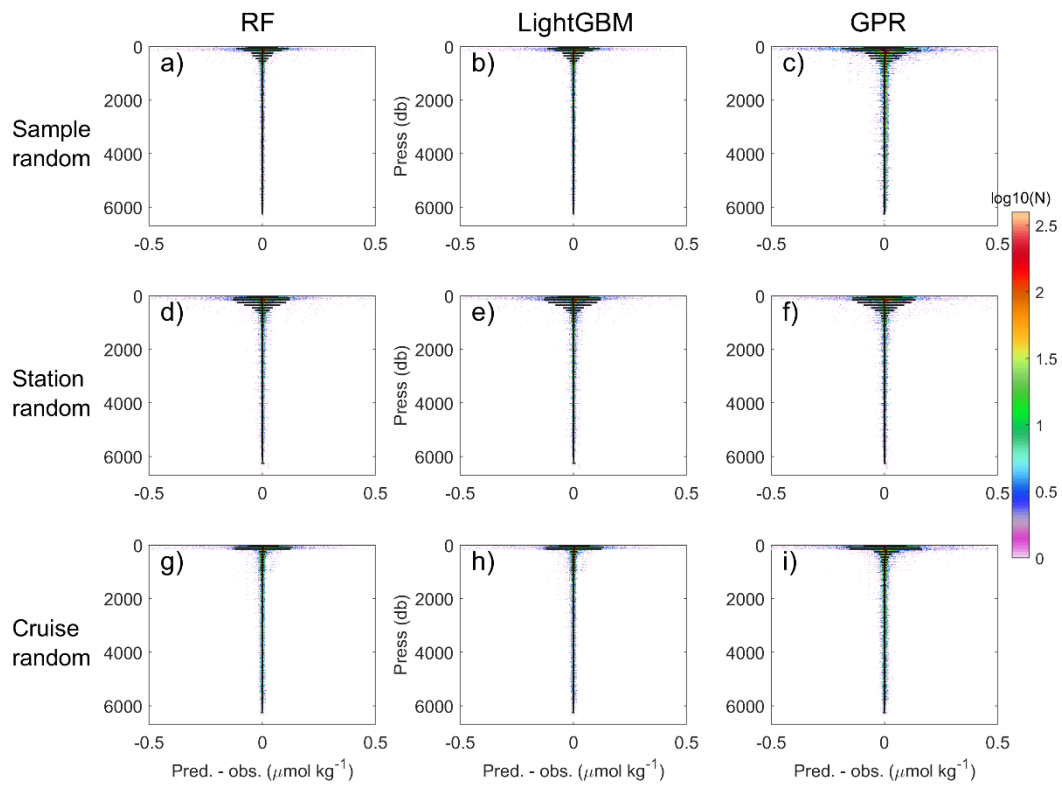
28



29

30 **Figure S4.** Vertical profiles of NO_x^- ($\text{NO}_3^- + \text{NO}_2^-$) reconstruction errors with different
 31 data selection strategies and machine learning methods. Colors indicate data density;
 32 the solid black vertical line represents the mean absolute error ($\mu\text{mol kg}^{-1}$); and the
 33 horizontal error bars denote the ± 1 standard deviation range. Plots in row 1 correspond
 34 to the sample random strategy (a-c); plots in row 2 correspond to the station random
 35 strategy (d-f); and plots in row 3 correspond to the cruise random strategy (g-i). Plots
 36 in column 1 correspond to the Random Forest (RF; a, d, g); plots in column 2
 37 correspond to LightGBM (b, e, h); and plots in column 3 correspond to Gaussian
 38 Process Regression (GPR; c, f, i).

39

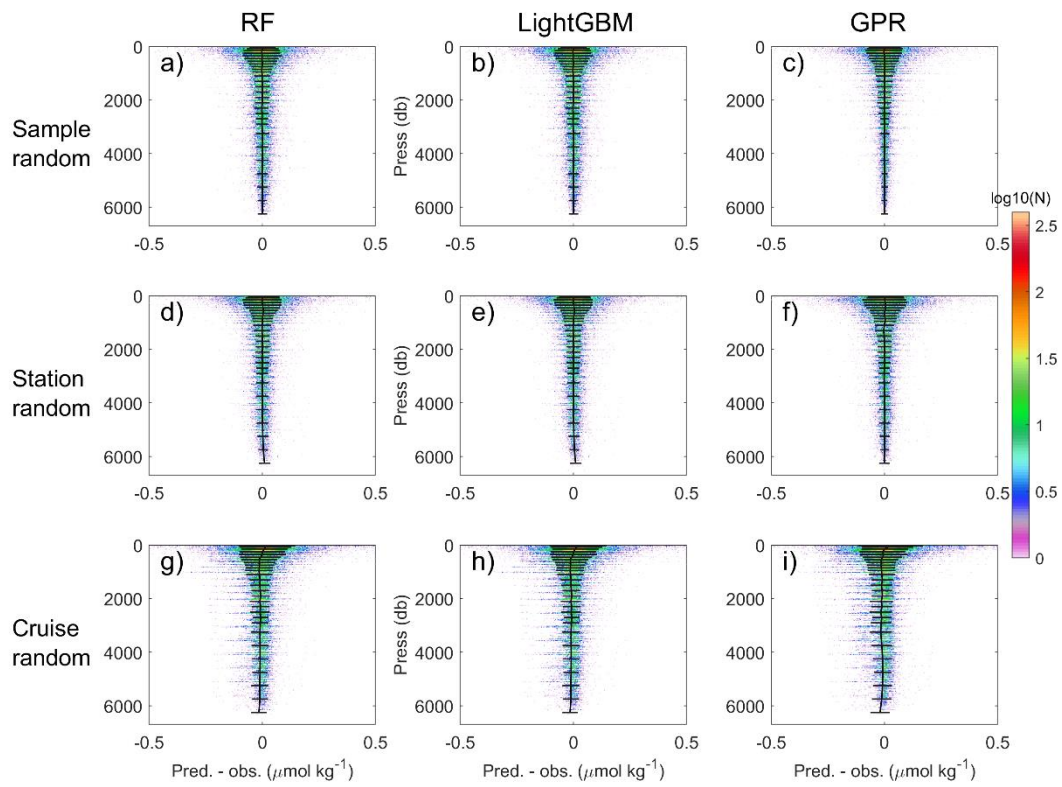


40

41 **Figure S5.** Same as Fig. S4, but for NO_2^- .

42

43

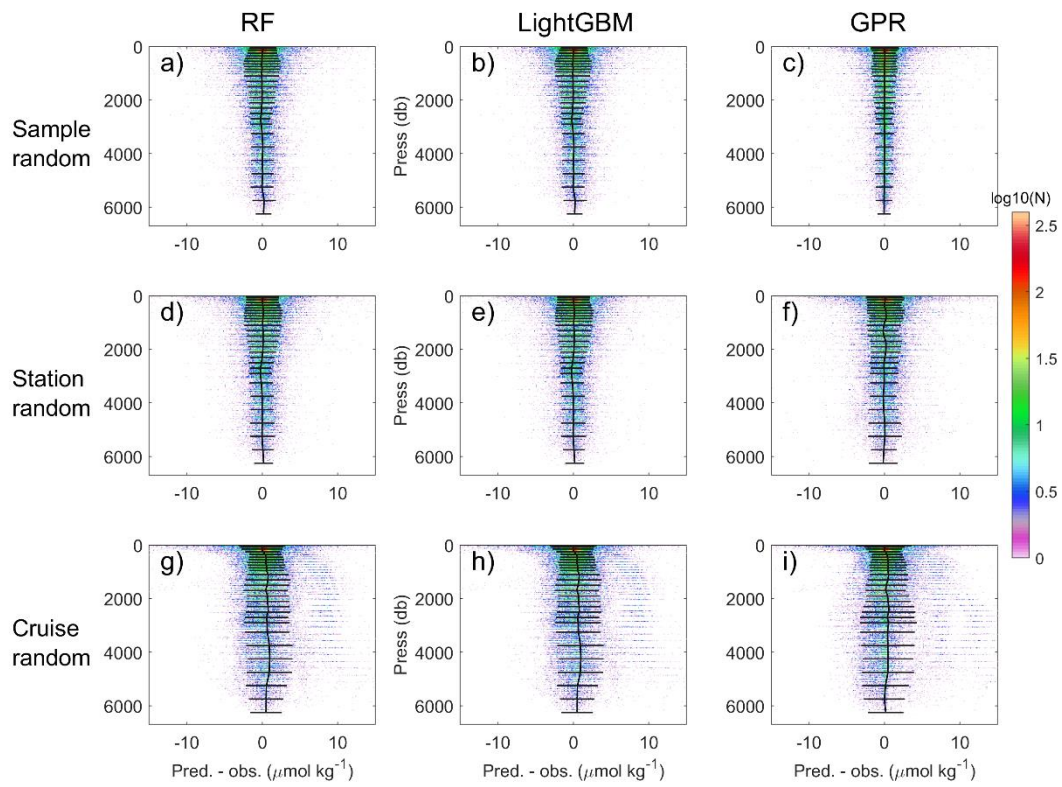


44

45 **Figure S6.** Same as Fig. S4, but for DIP.

46

47

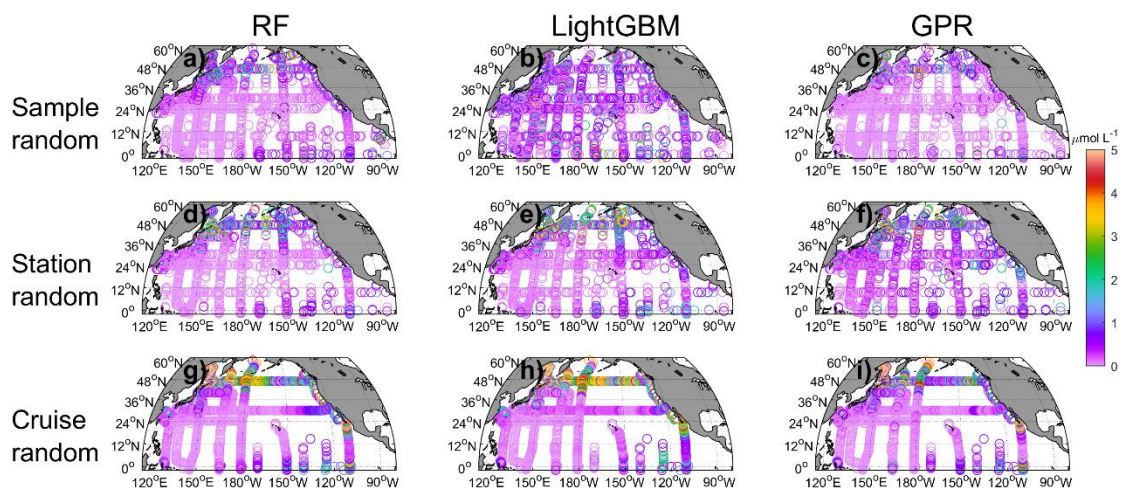


48

49 **Figure S7.** Same as Fig. S4, but for Si(OH)_4 .

50

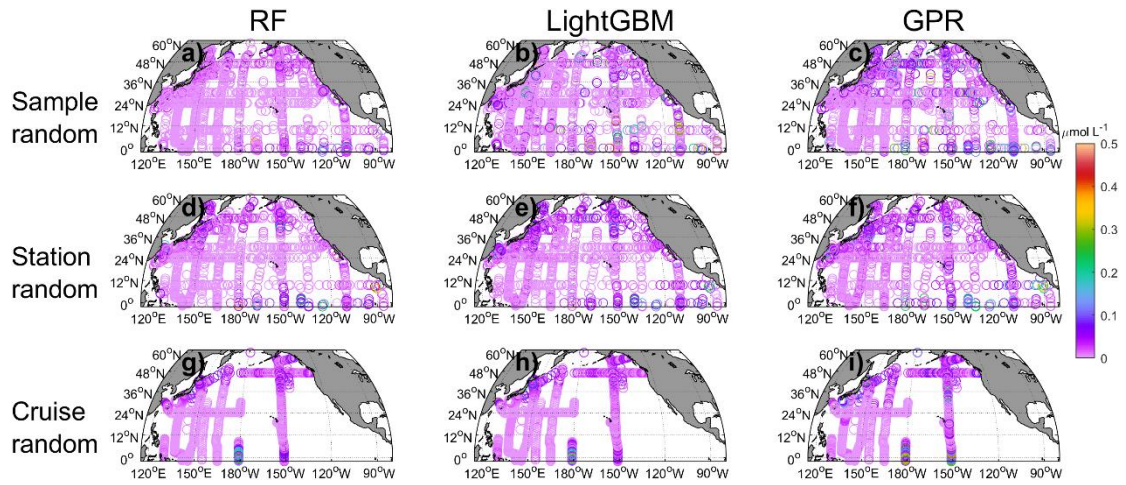
51



52

53 **Figure S8.** Spatial distribution of absolute errors (absolute value of predicted subtract
 54 observed values) for surface NO_x^- concentrations with different data selection strategies
 55 and machine learning methods. Plots shown in row 1 correspond to the sample random
 56 strategy (a-c), row 2 correspond to the station random strategy (d-e), and row 3
 57 correspond to the cruise random strategy (g-i). Plots shown in column 1 correspond to
 58 the Random Forest (RF; a, d, and g), column 2 correspond to the LightGBM (b, e, and
 59 h), and column 3 correspond to the Gaussian Process Regression (GPR; c, f, and i).

60

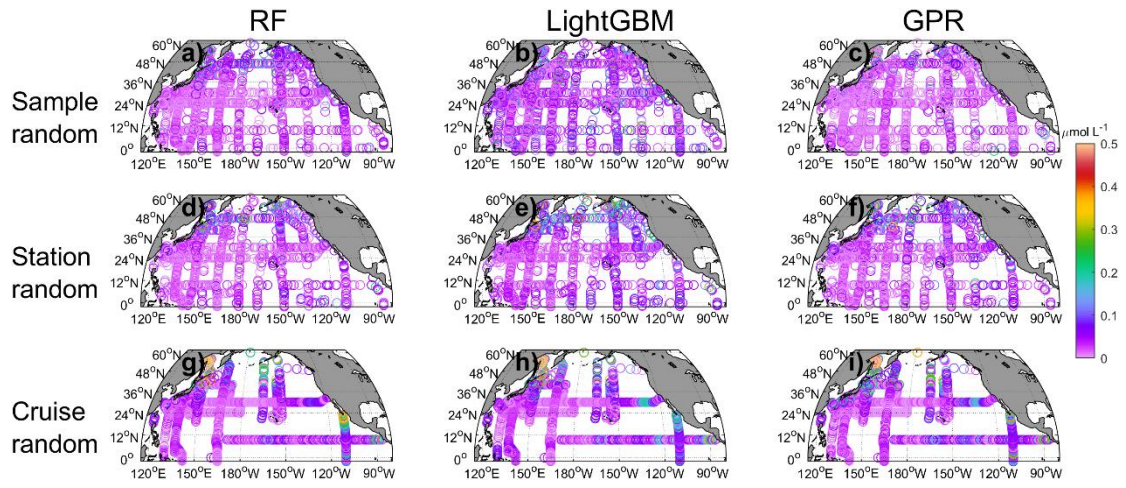


61

62 **Figure S9.** Same as Fig. S8, but for NO_2^- .

63

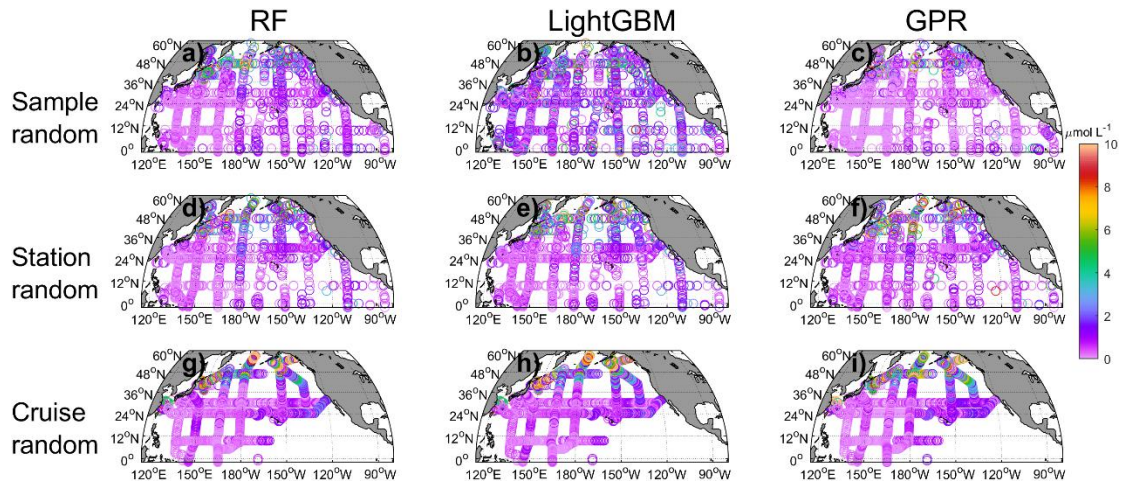
64



65

66 **Figure S10.** Same as Fig. S8, but for DIP.

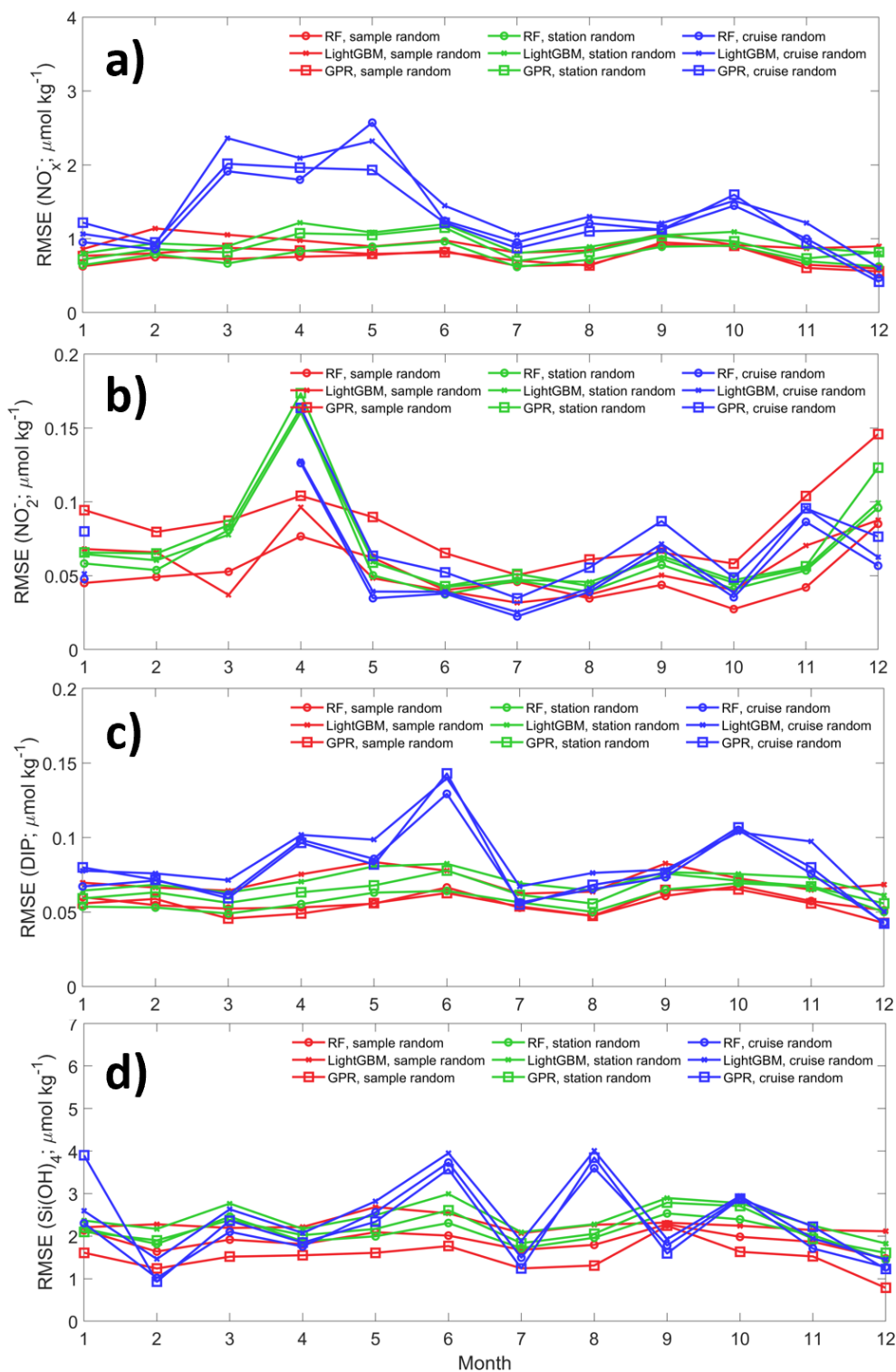
67



68

69 **Figure S11.** Same as Fig. S8, but for Si(OH)_4 .

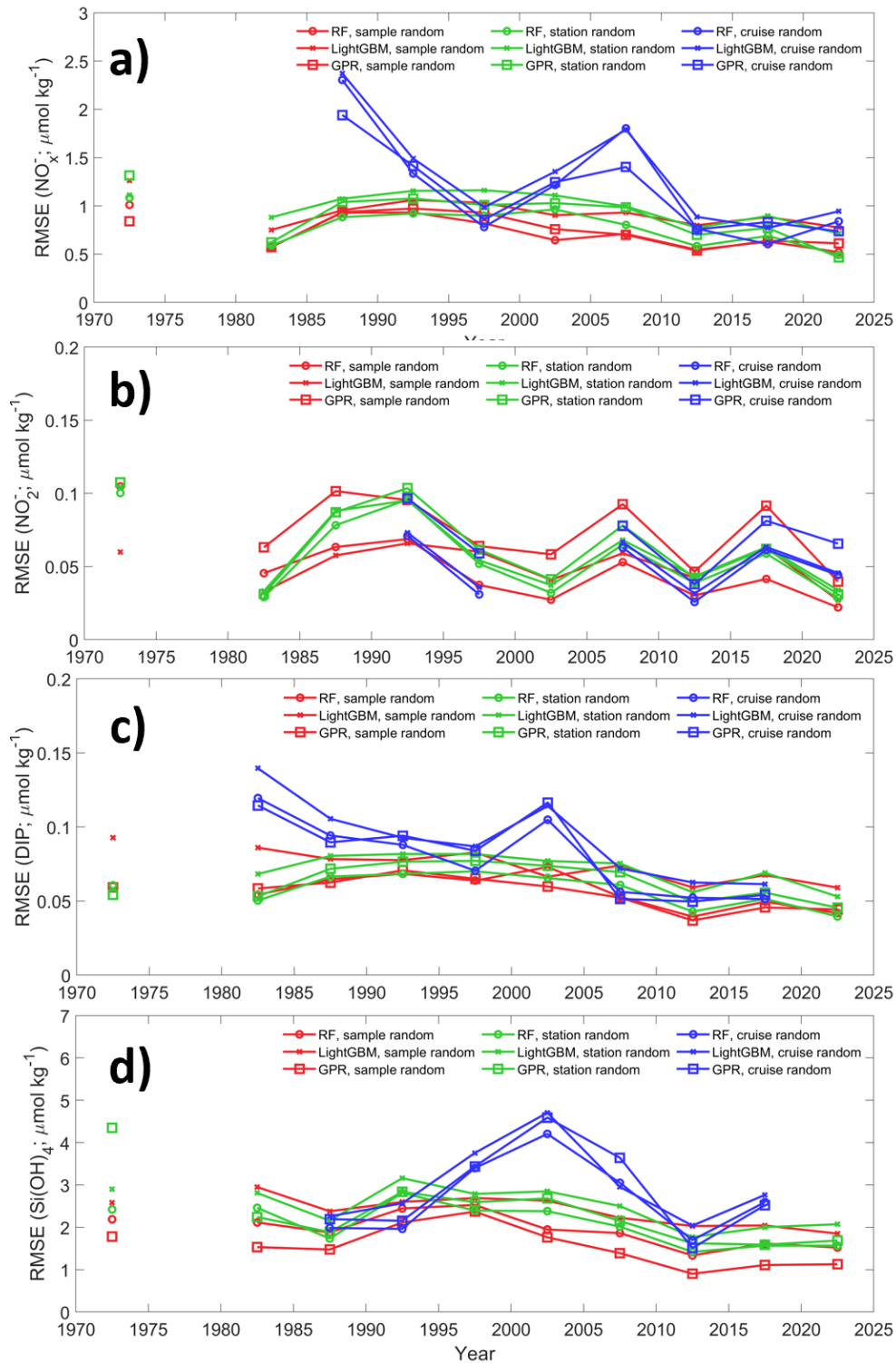
70



71

72 **Figure S12.** Monthly variations in reconstruction errors (root mean square error) for
 73 different parameters under various data selection strategies and machine learning
 74 methods.

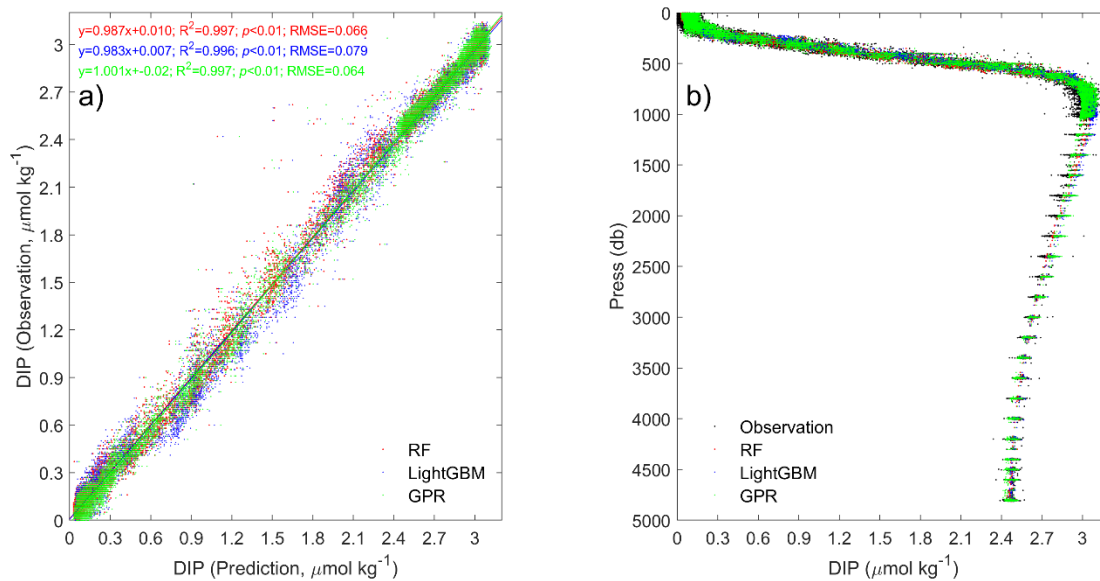
75



76

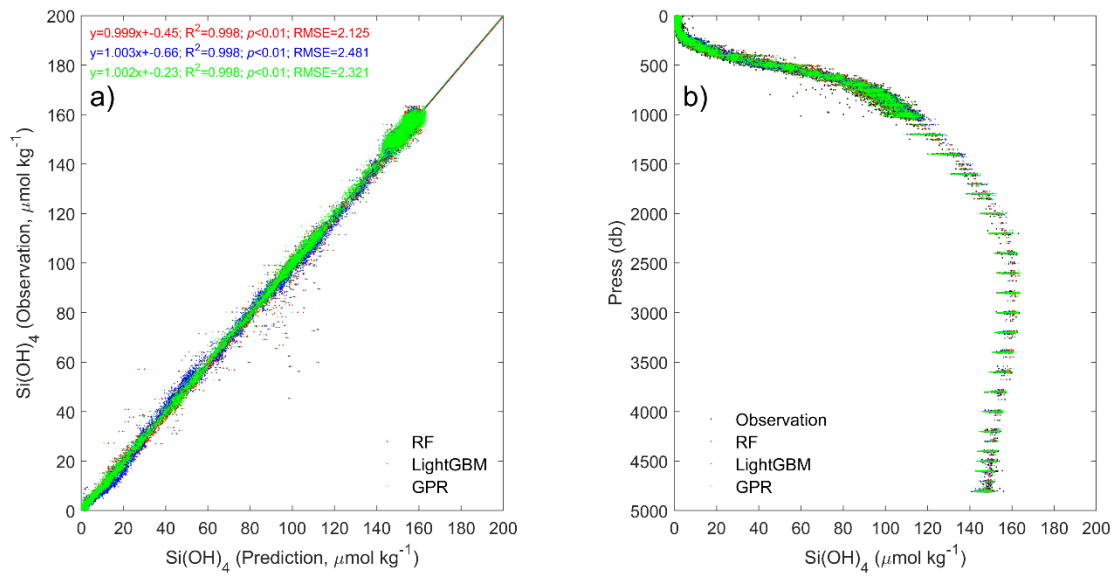
77 **Figure S13.** Long-term variations in reconstruction error (root mean square error)
 78 under different data selection strategies and machine learning methods.

79



80
 81 **Figure S14.** Validating the reconstructed DIP concentrations at Station ALOHA. a)
 82 Reconstructed DIP vs. observations: Random Forest (RF; red dots), LightGBM (blue
 83 dots), and Gaussian Process Regression (GPR; green dots); b) Profiles of observed
 84 (black dots) and reconstructed DIP from RF (red dots), LightGBM (blue dots), and GPR
 85 (green dots).

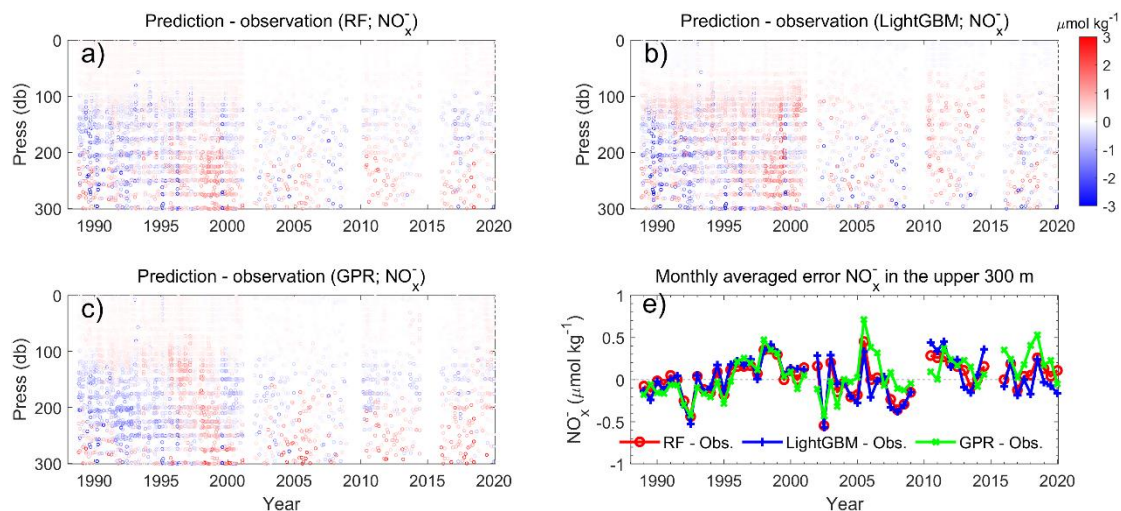
86
 87
 88



89

90 **Figure S15.** Similar to Fig. S14, but for Si(OH)_4 .

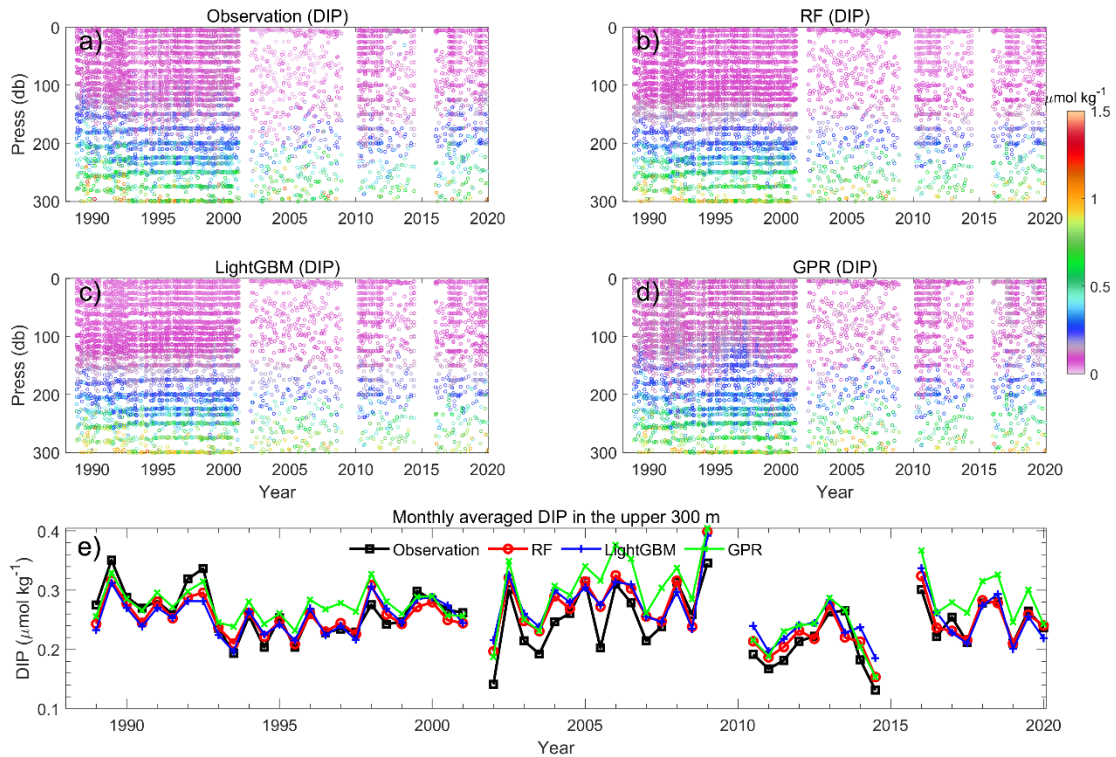
91



92

93 **Figure S16.** Temporal variations of reconstruction error profiles of NO_x^- ($\text{NO}_3^- + \text{NO}_2^-$)
 94 (prediction minus observations) in the upper 300 m at Station ALOHA from 1988 to
 95 2021 with different machine learning models: Random Forest (RF; a), LightGBM (b),
 96 and Gaussian Process Regression (GPR; c). Note that the colorbar plotted in (b) is also
 97 applicable to (a) and (c). (e) Temporal variations of monthly averaged NO_x^- prediction
 98 errors (prediction minus observations) in the upper 300 m from RF, LightGBM, and
 99 GPR.

100

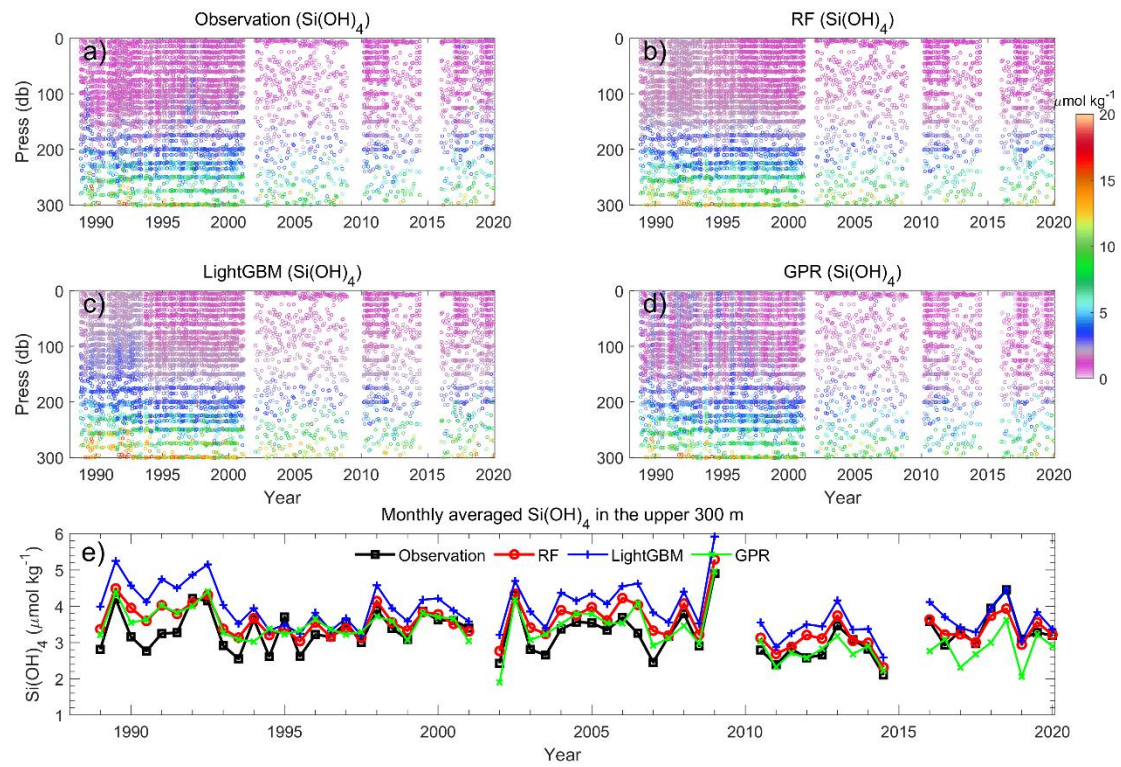


101

102 **Figure S17.** Temporal variations of DIP concentrations in the upper 300 m at Station
 103 ALOHA from 1988 to 2021 for observed (a) and reconstructed DIP by Random Forest
 104 (RF; b), LightGBM (c), and Gaussian Process Regression (GPR; d). (e) Time series of
 105 monthly averaged NO_x^- concentrations in the upper 300 m from observations, and
 106 reconstructions by RF, LightGBM, and GPR, respectively.

107

108



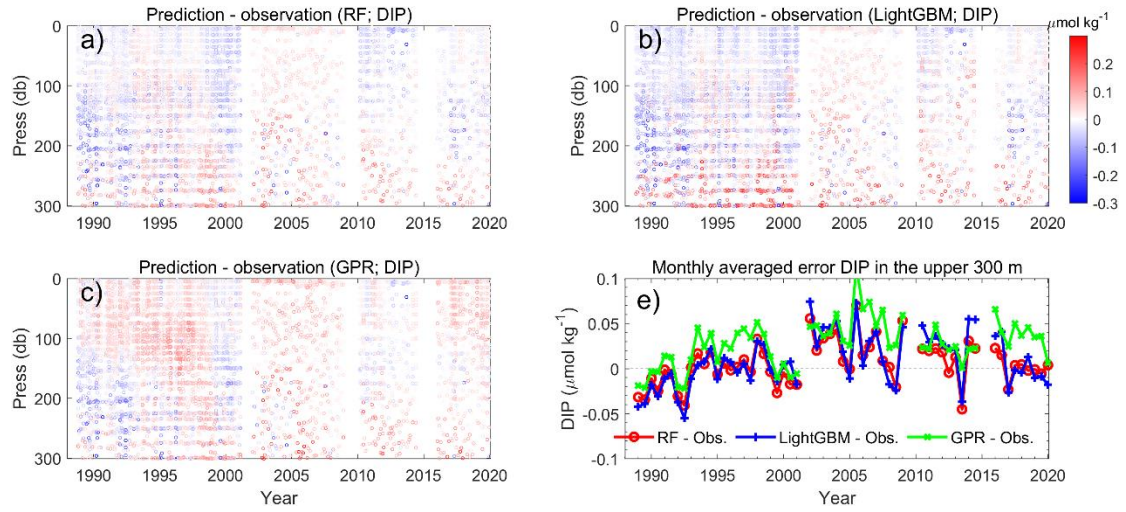
109

110 **Figure S18.** Similar to Fig. S17, but for Si(OH)_4 .

111

112

113

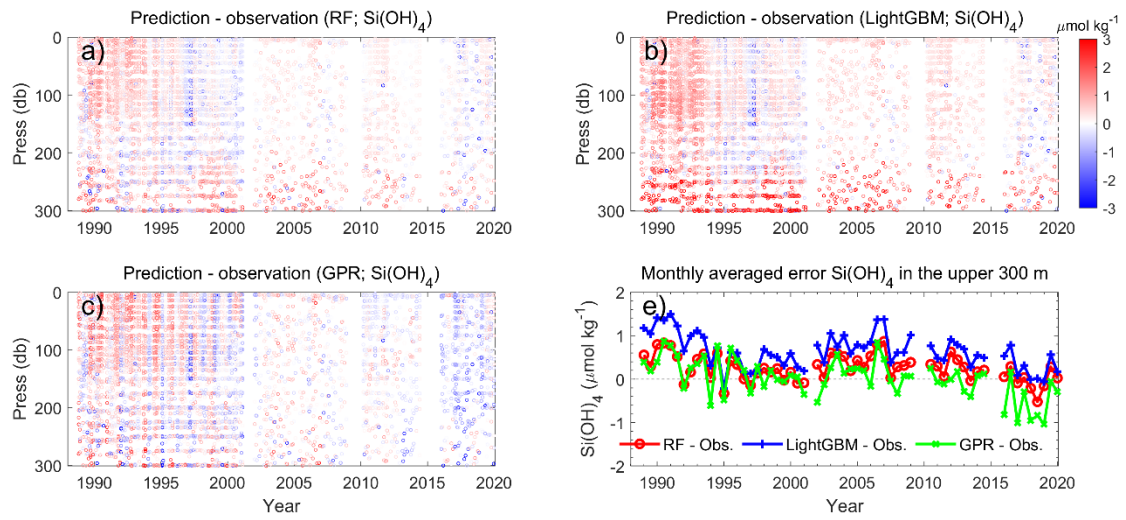


114

115 **Figure S19.** Same as Fig. S16, but for DIP.

116

117

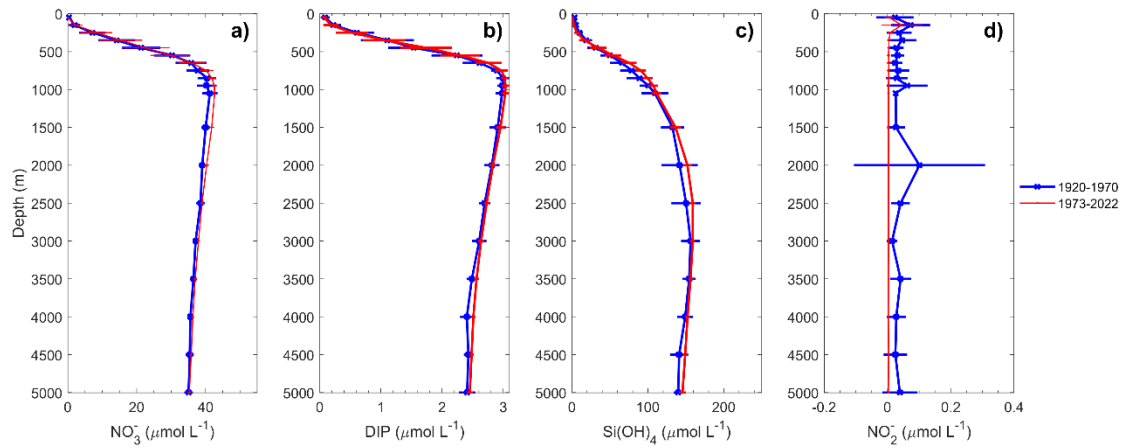


118

119 **Figure S20.** Similar to Fig. S16, but for Si(OH)_4 .

120

121



122

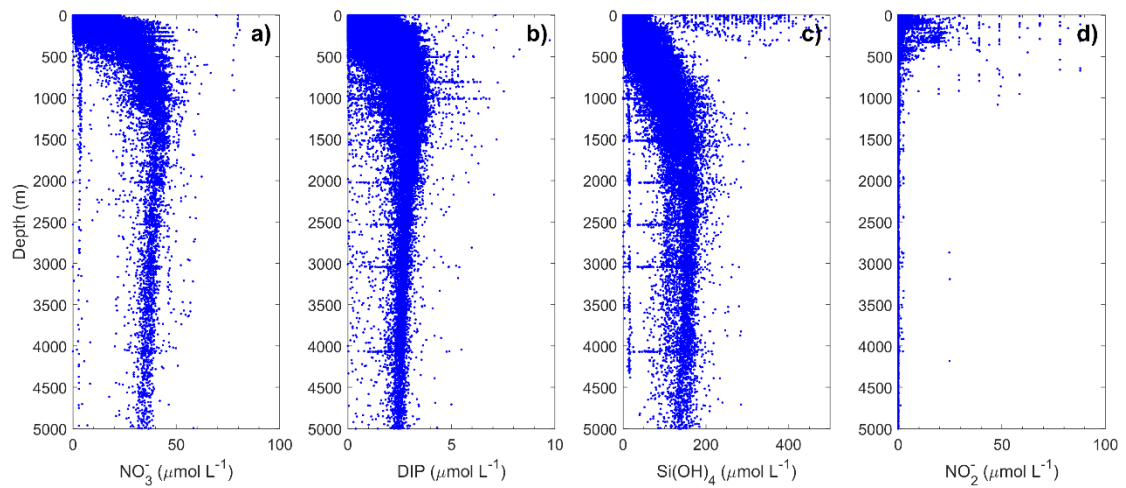
123 **Figure S21.** Mean profiles (vertical lines) with standard deviations (horizontal lines) of

124 nutrients between 1920–1970 (Ocean Station Data in the World Ocean Database after

125 quality control) and 1973–2022 (CCHDO data) in the North Pacific basin region

126 (180°E – 150°W , 15°N – 30°N): (a) NO_3^- ; (b) DIP; (c) Si(OH)_4 ; (d) NO_2^- .

127



128

129

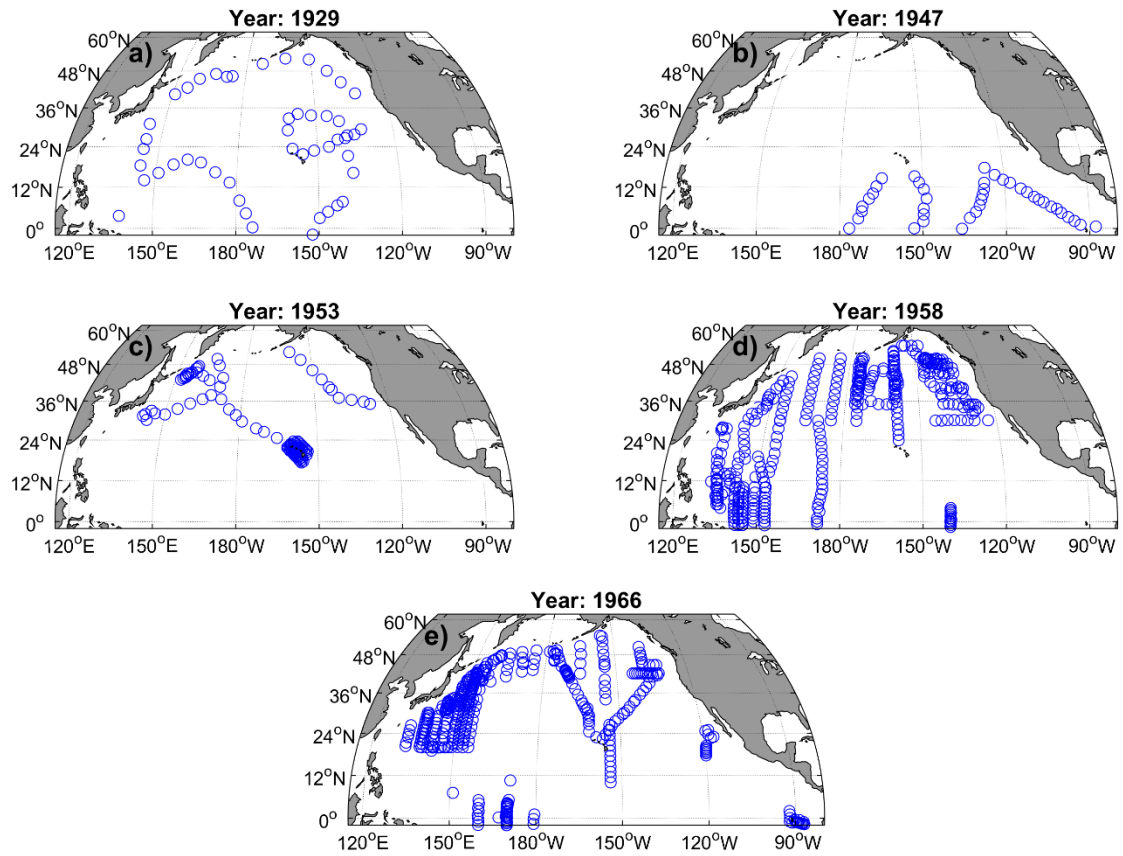
130

131

132

133

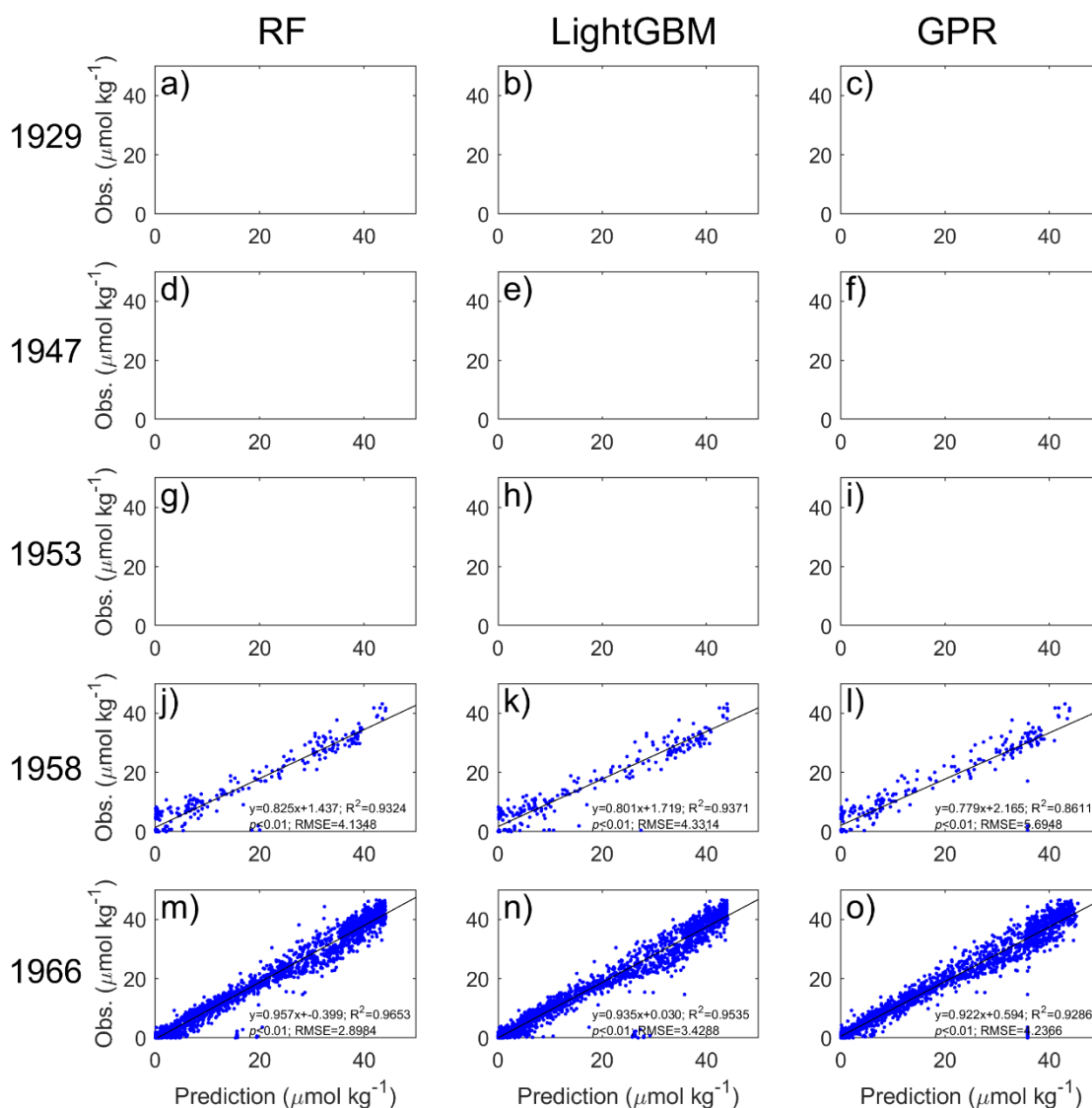
Figure S22. Profiles of nutrients collected from the Ocean Station Data in the World Ocean Database (before 1970) in the North Pacific basin: (a) NO_3^- ; (b) DIP; (c) Si(OH)_4 ; (d) NO_2^- , illustrating the scattered vertical distribution.



134

135 **Figure S23.** Hydrographic and nutrient sampling locations for selected years prior to
 136 1970 (collected from Ocean Station Data in the World Ocean Database): (a) 1929, (b)
 137 1947, (c) 1953, (d) 1958, and (e) 1966.

138

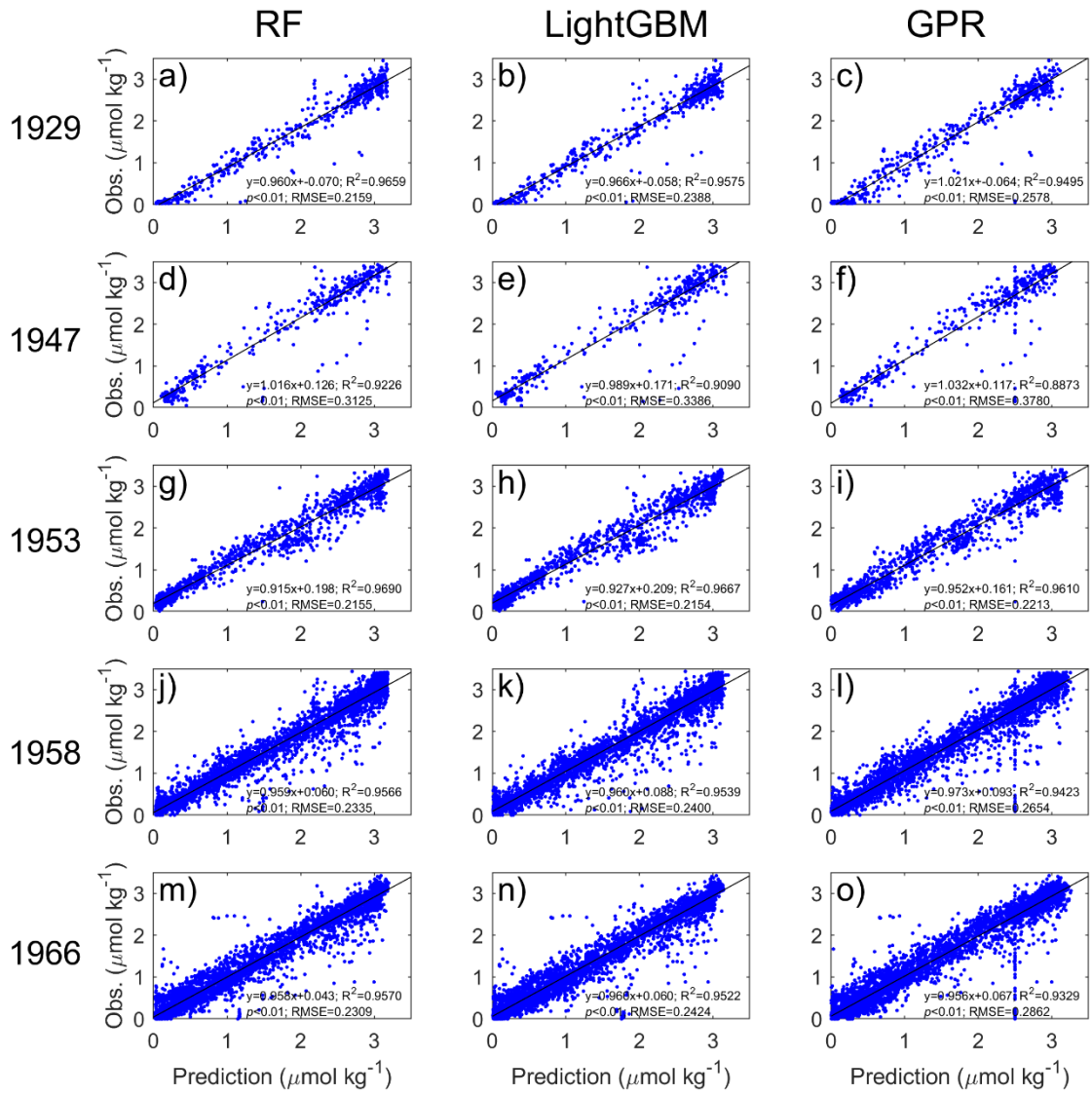


139

140 **Figure S24.** Validation of reconstructed NO_3^- concentrations from hydrography-based
 141 reconstructions for 1929 (row 1), 1947 (row 2), 1953 (row 3), 1958 (row 4), and 1966
 142 (row 5). Station locations are shown in Fig. S23. Columns correspond to different
 143 reconstruction methods: Random Forest (RF; a, d, g, j, m), LightGBM (b, e, h, k, n),
 144 and Gaussian Process Regression (GPR; c, f, i, l, o). Each panel displays the linear
 145 regression fit (black line), regression equation, coefficient of determination (R^2), p-
 146 value, and root mean square error (RMSE). Note: no NO_3^- observations were available
 147 for 1929, 1947, or 1953.

148

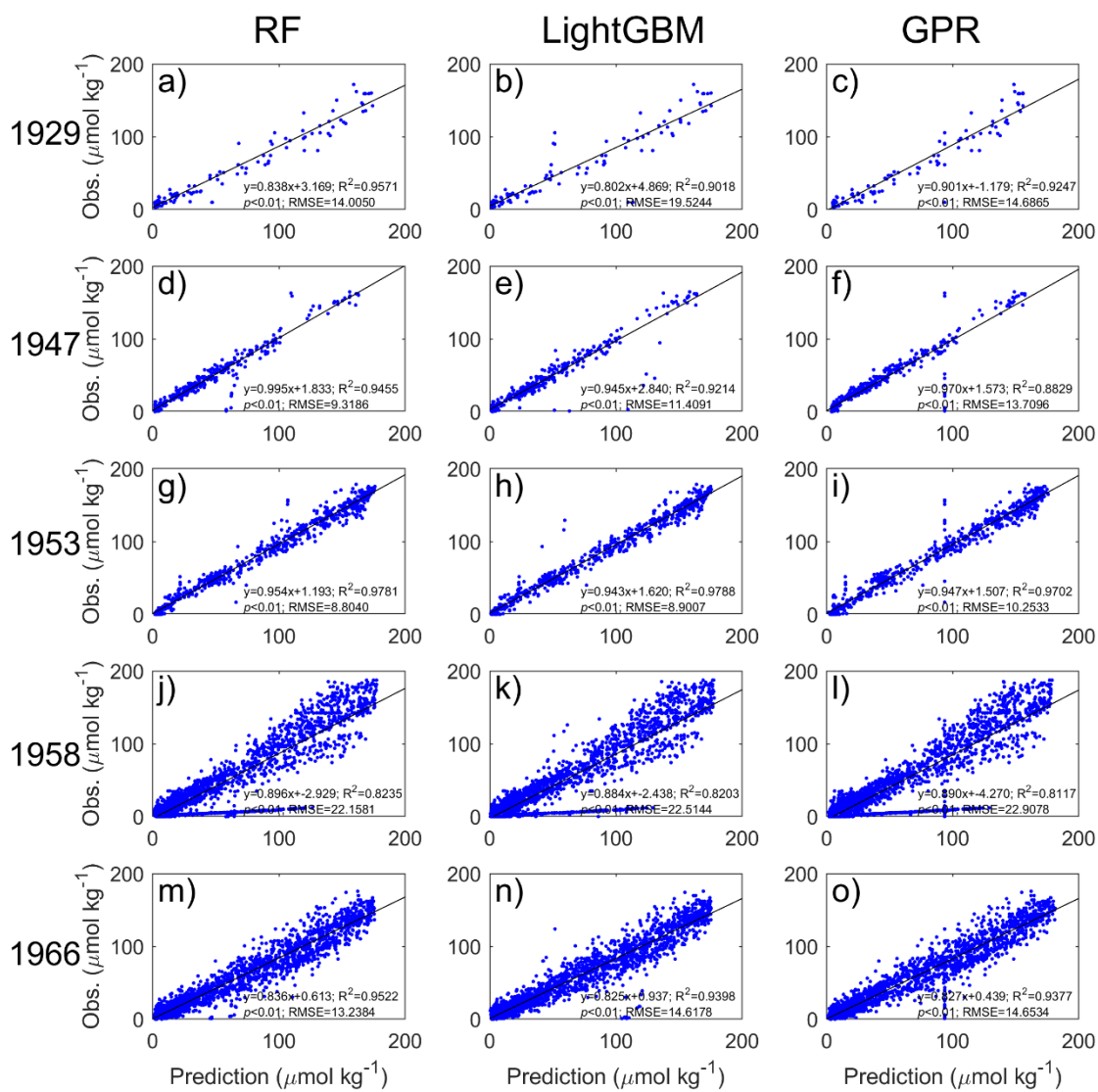
149



150

151 **Figure S25.** Same as Fig. S24, but for DIP.

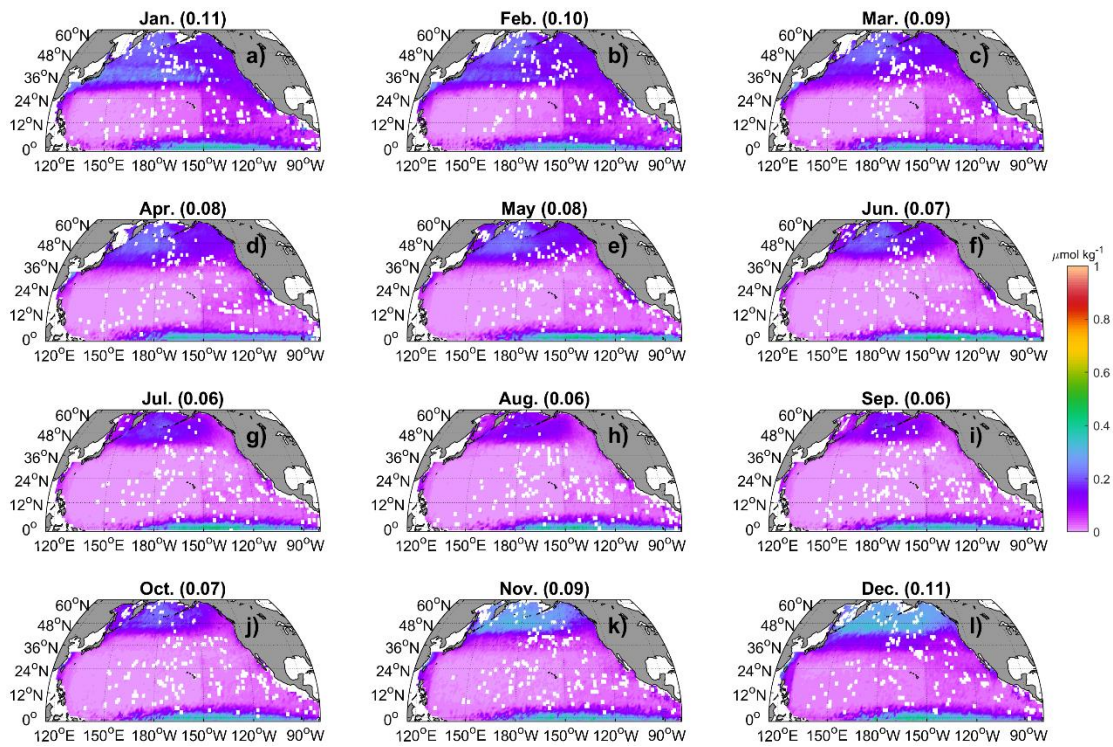
152



153

154 **Figure S26.** Same as Fig. S24, but for Si(OH)_4 .

155



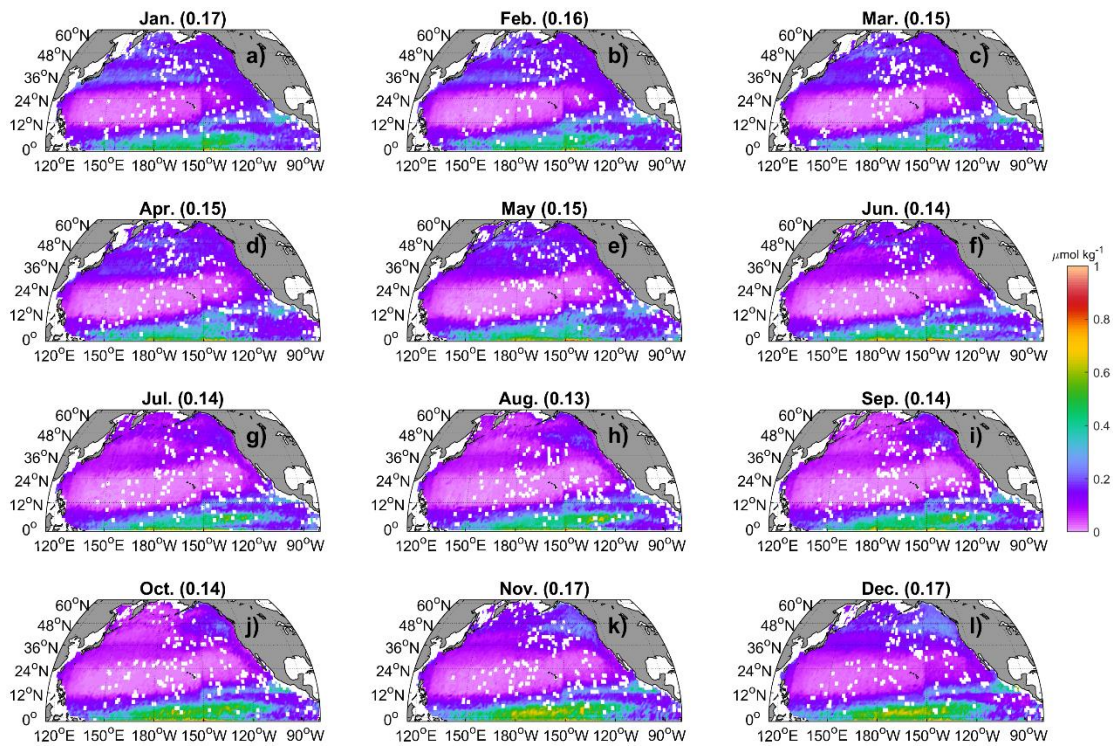
156

157 **Figure S27.** The monthly climatology of NO_2^- at 5 m in the North Pacific. Data are

158 binned and averaged within $1 \times 1^\circ$ grid cell. The values in the title represent the spatial

159 mean values.

160

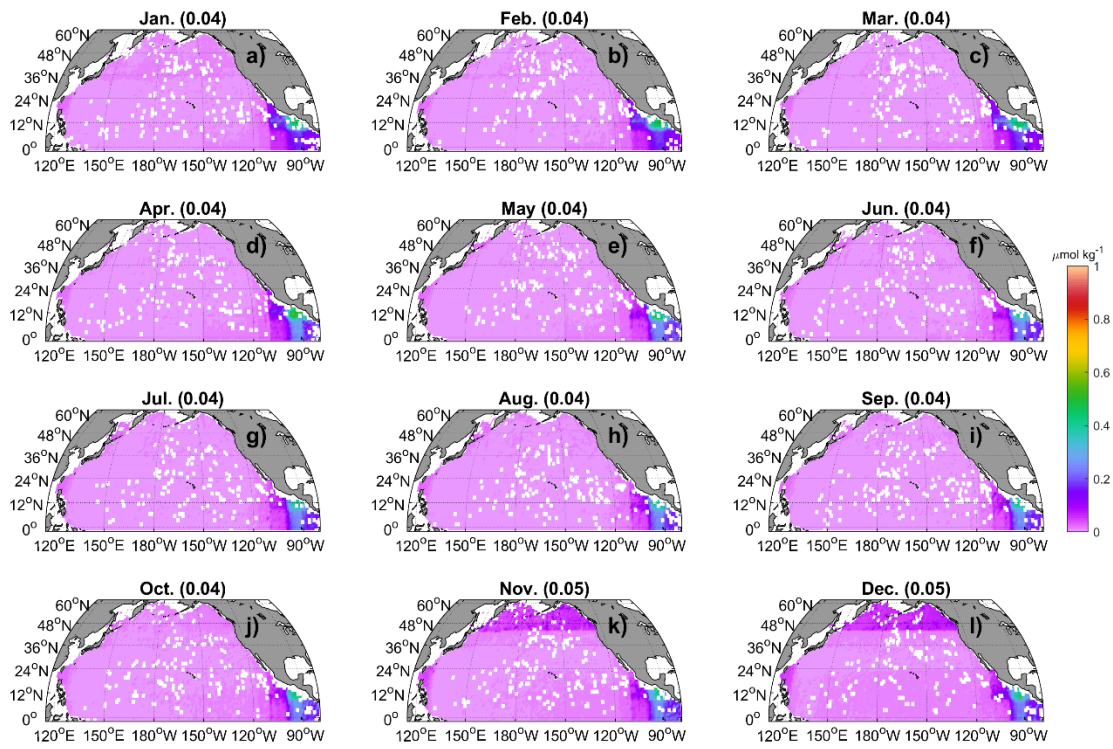


161

162 **Figure S28.** Similar to Fig. S27, but for a depth of 100 m.

163

164

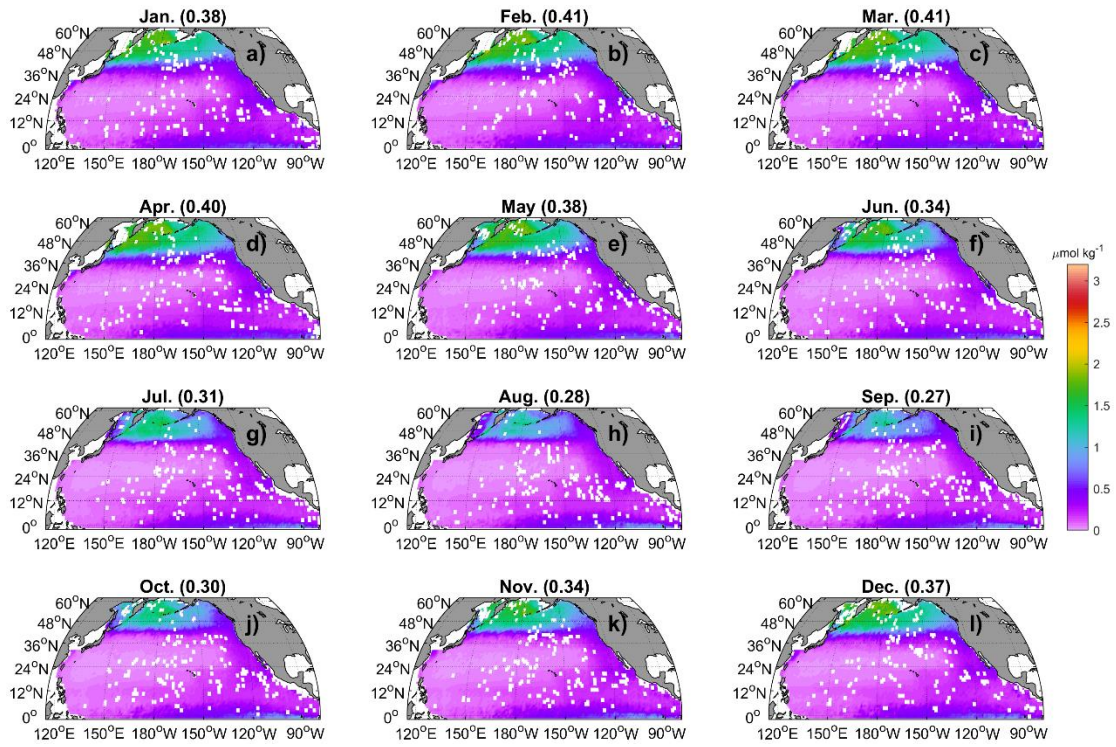


165

166 **Figure S29.** Similar to Fig. S27, but for a depth of 500 m.

167

168



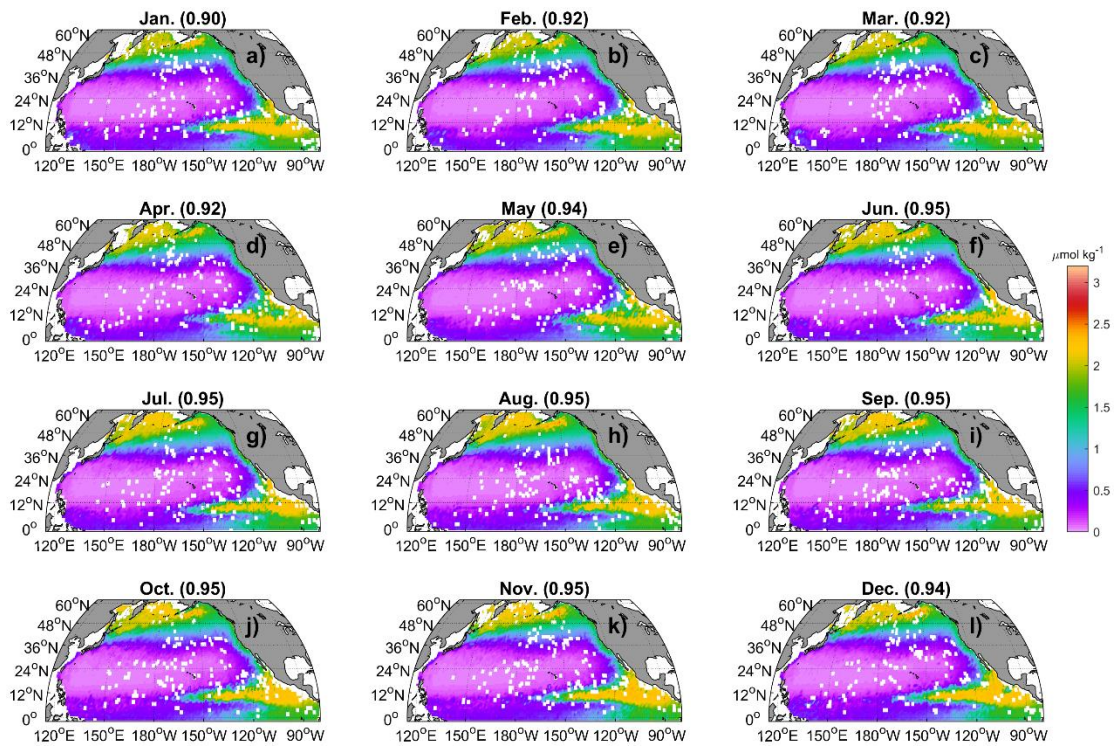
169

170

Figure S30. Similar to Fig. S27, but for DIP and at a depth of 5 m.

171

172



173

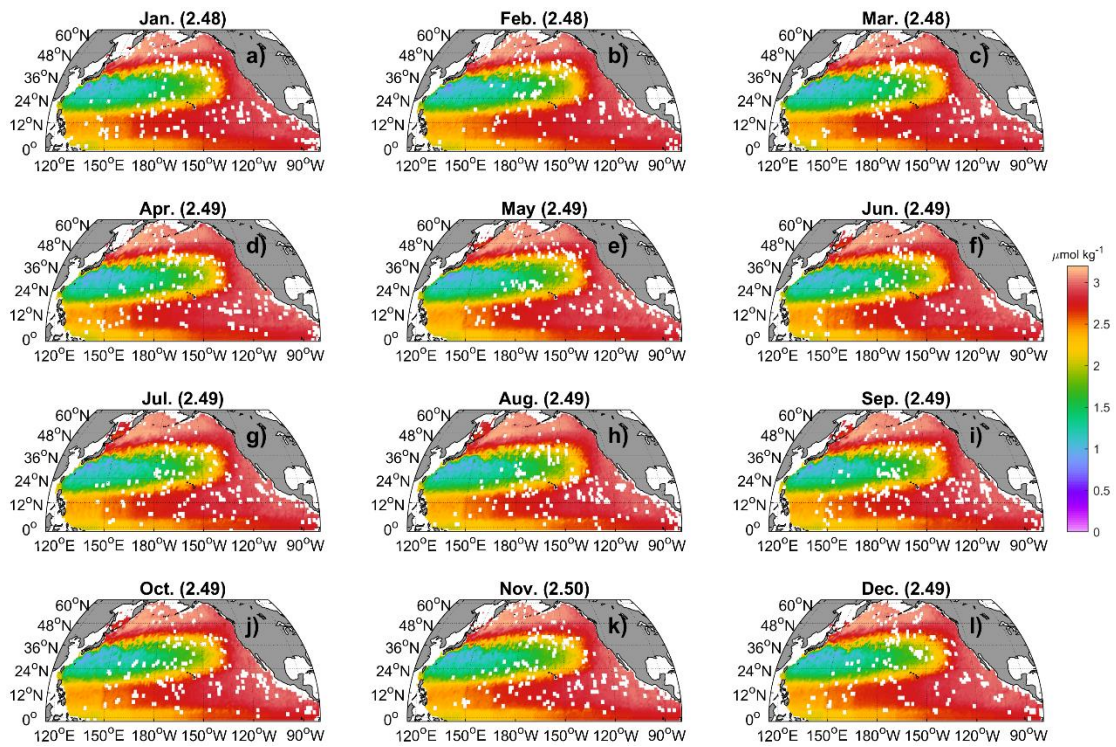
174 **Figure S31.** Similar to Fig. S27, but for DIP and at a depth of 100 m.

175

176

177

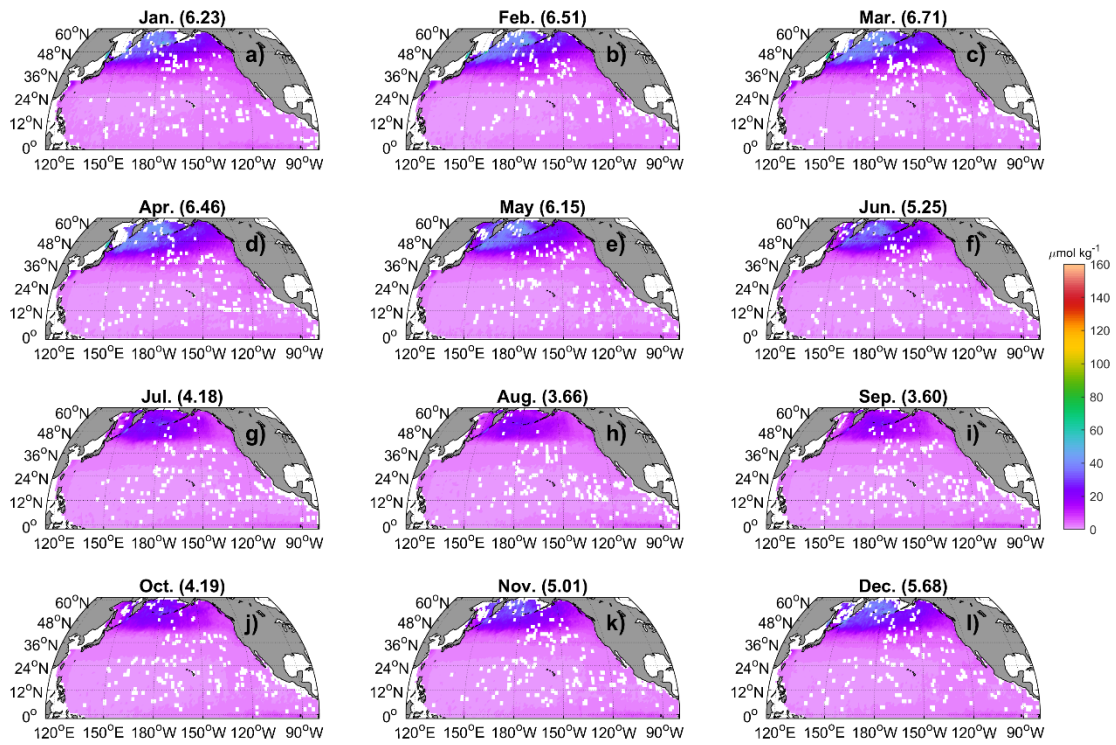
178



179

180 **Figure S32.** Similar to Fig. S27, but for DIP and at a depth of 500 m.

181

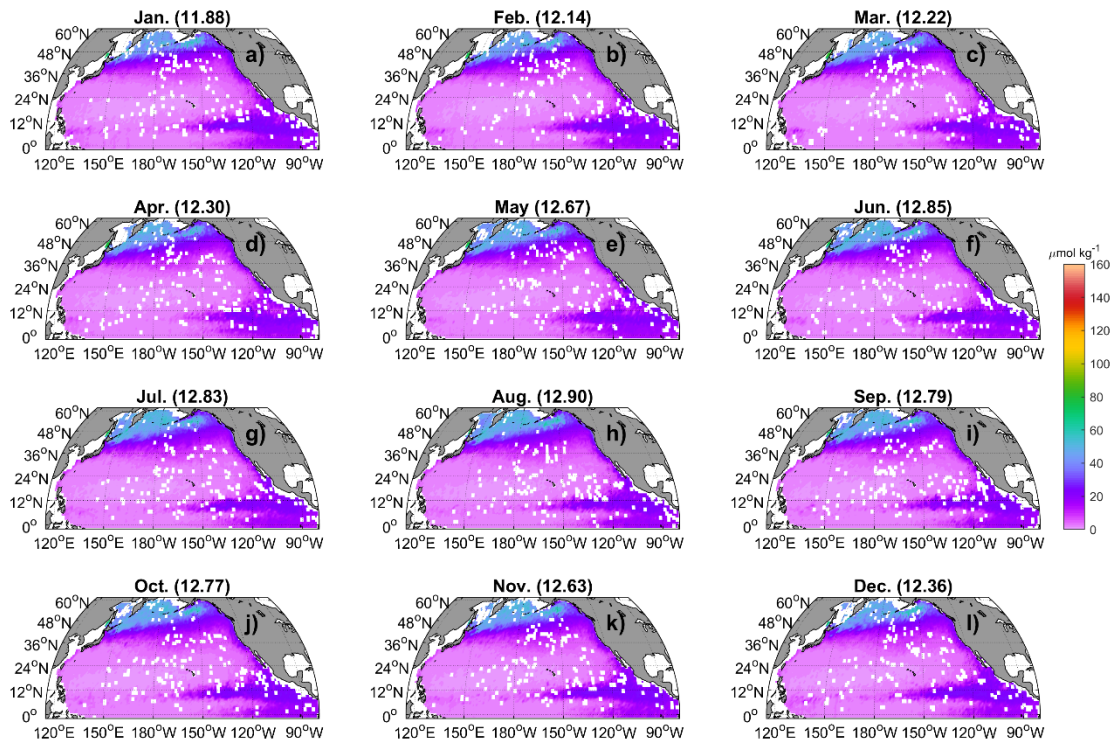


182

183 **Figure S33.** Similar to Fig. S27, but for Si(OH)_4 and at a depth of 5 m.

184

185



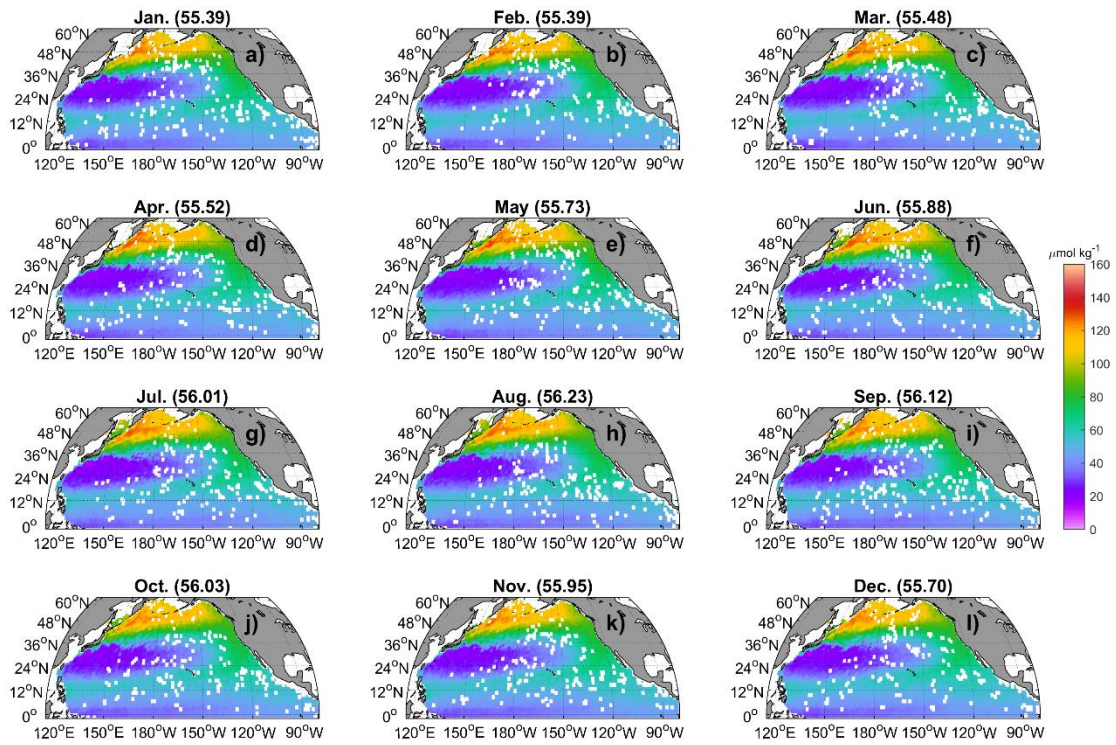
186

187 **Figure S34.** Similar to Fig. S27, but for Si(OH)_4 and at a depth of 100 m.

188

189

190

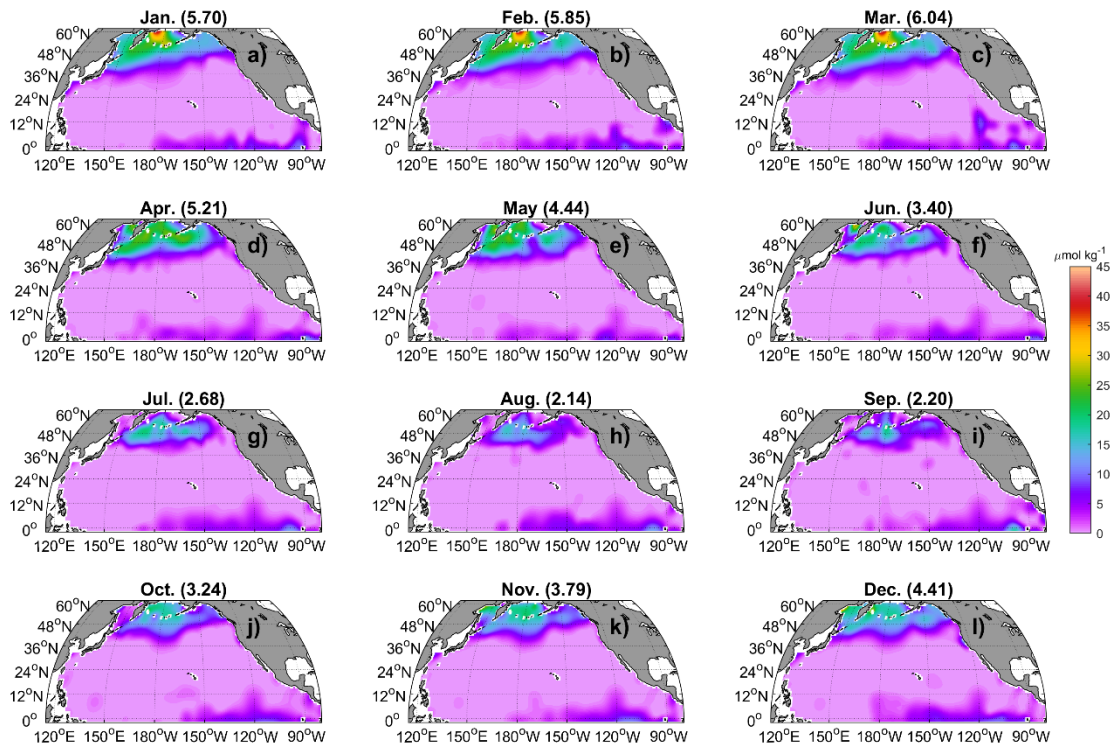


191

192 **Figure S35.** Similar to Fig. S27, but for Si(OH)_4 and at a depth of 500 m.

193

194



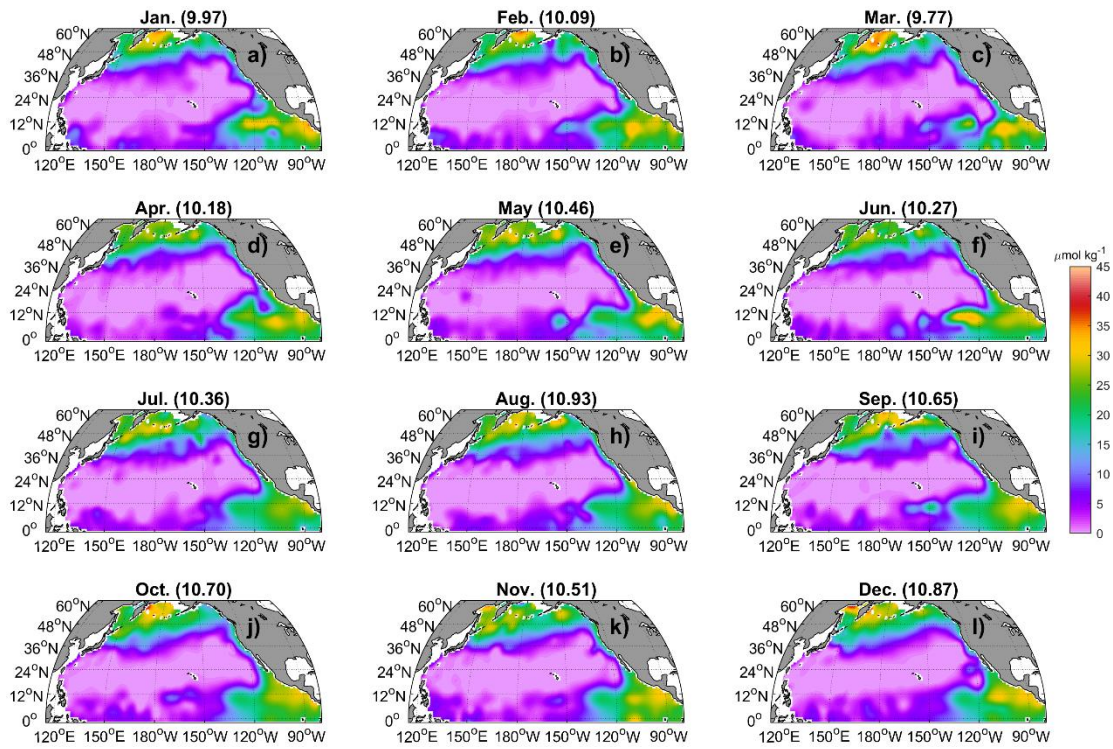
195

196 **Figure S36.** The WOA23 monthly climatology of NO_3^- at 5 m in the North Pacific.

197 The values in the title represent the spatial mean values.

198

199

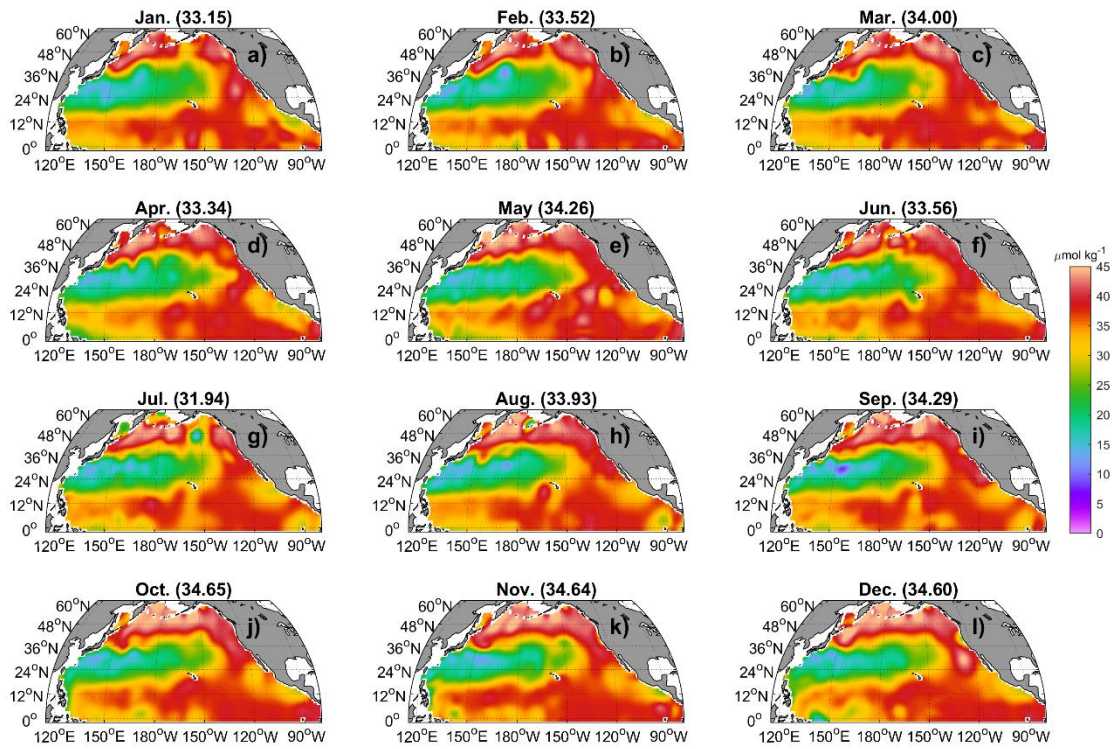


200

201 **Figure S37.** Similar to Fig. S36, but for a depth of 100 m.

202

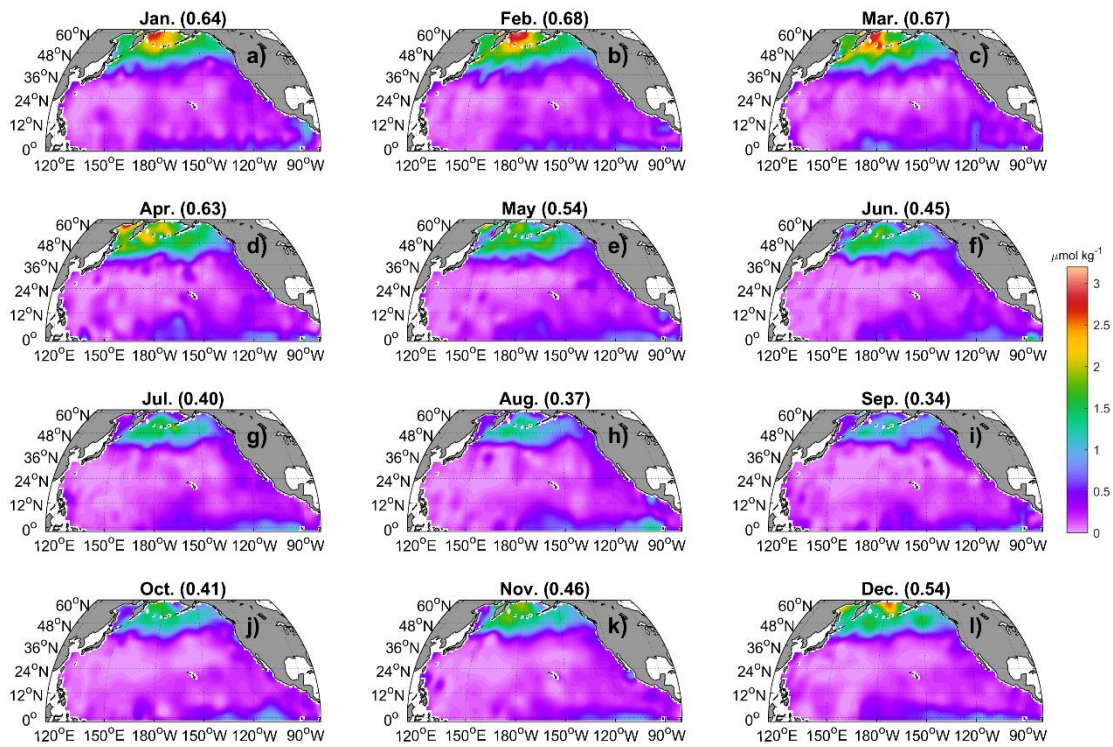
203



204

205 **Figure S38.** Similar to Fig. S36, but for a depth of 500 m.

206

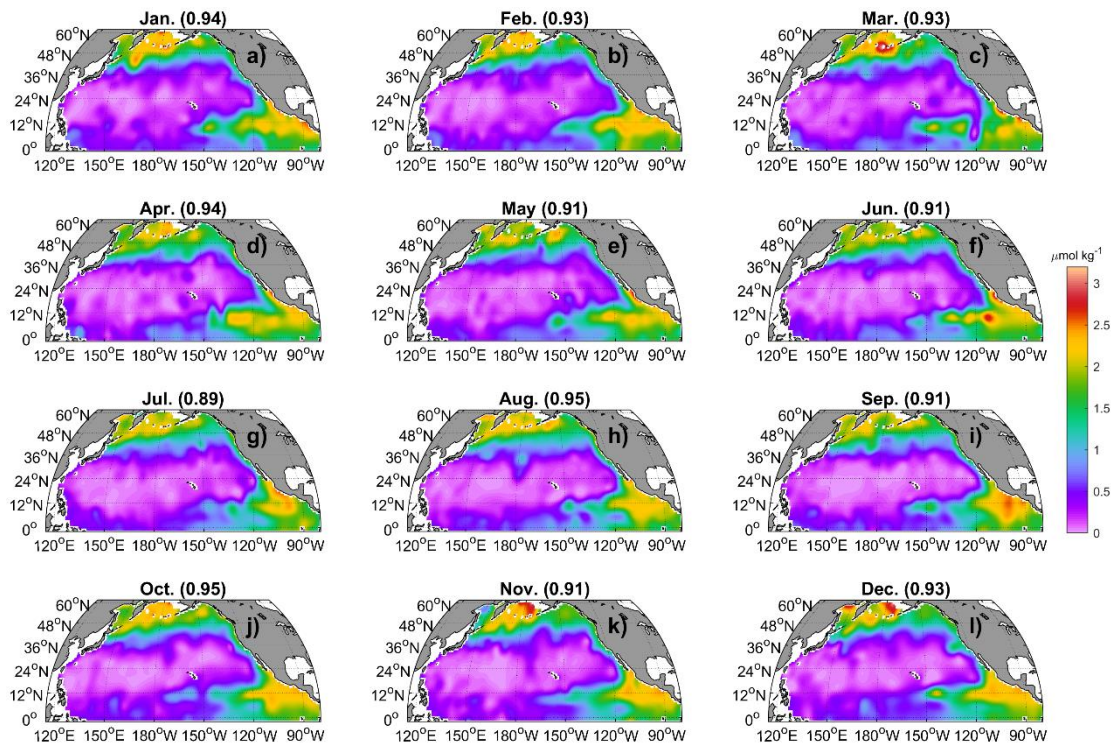


207

208 **Figure S39.** Similar to Fig. S36, but for DIP and at a depth of 5 m.

209

210

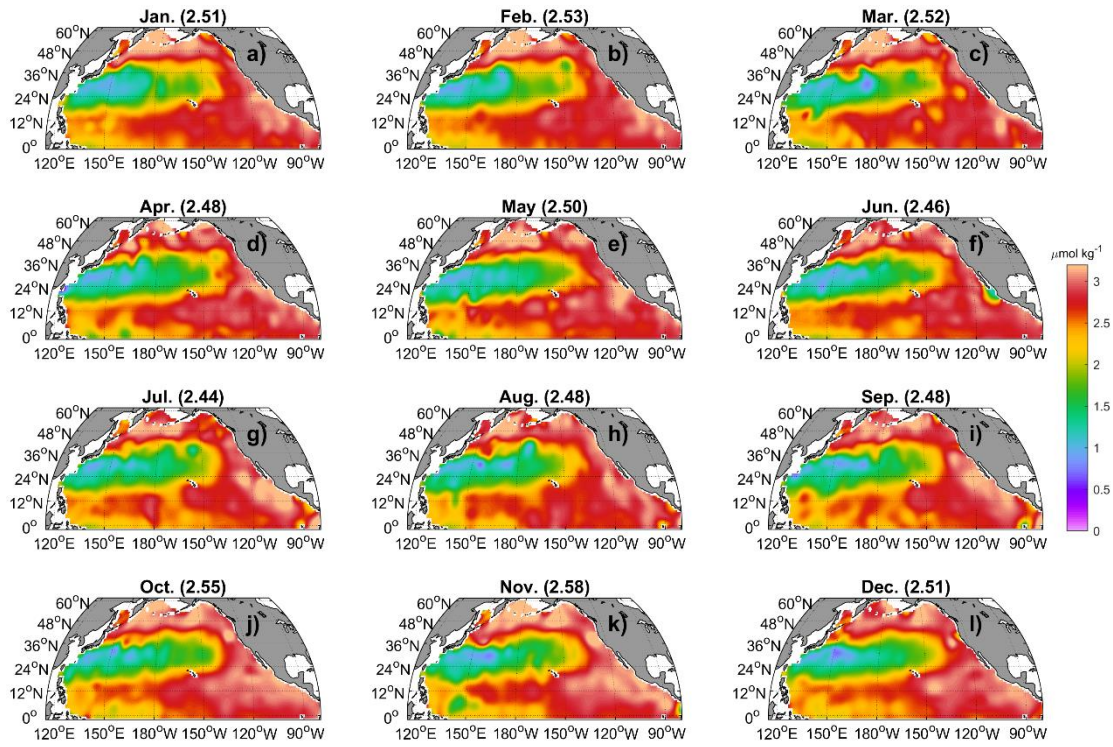


211

212 **Figure S40.** Similar to Fig. S36, but for DIP and at a depth of 100 m.

213

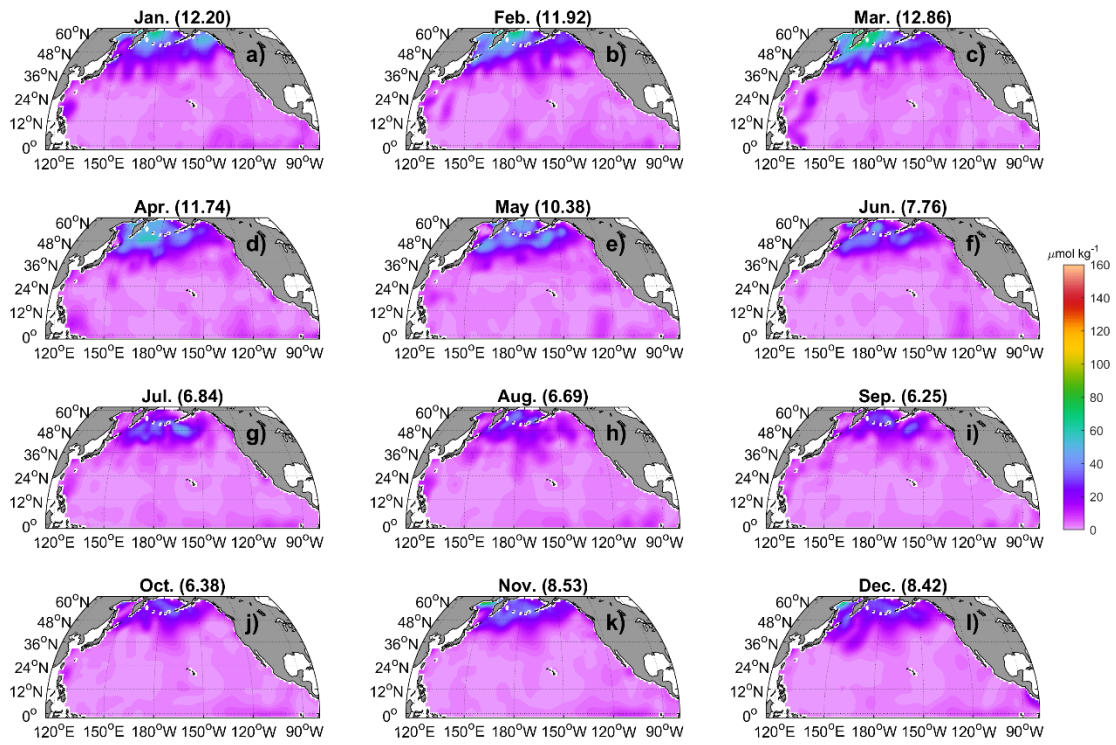
214



215

216 **Figure S41.** Similar to Fig. S36, but for DIP and at a depth of 500 m.

217

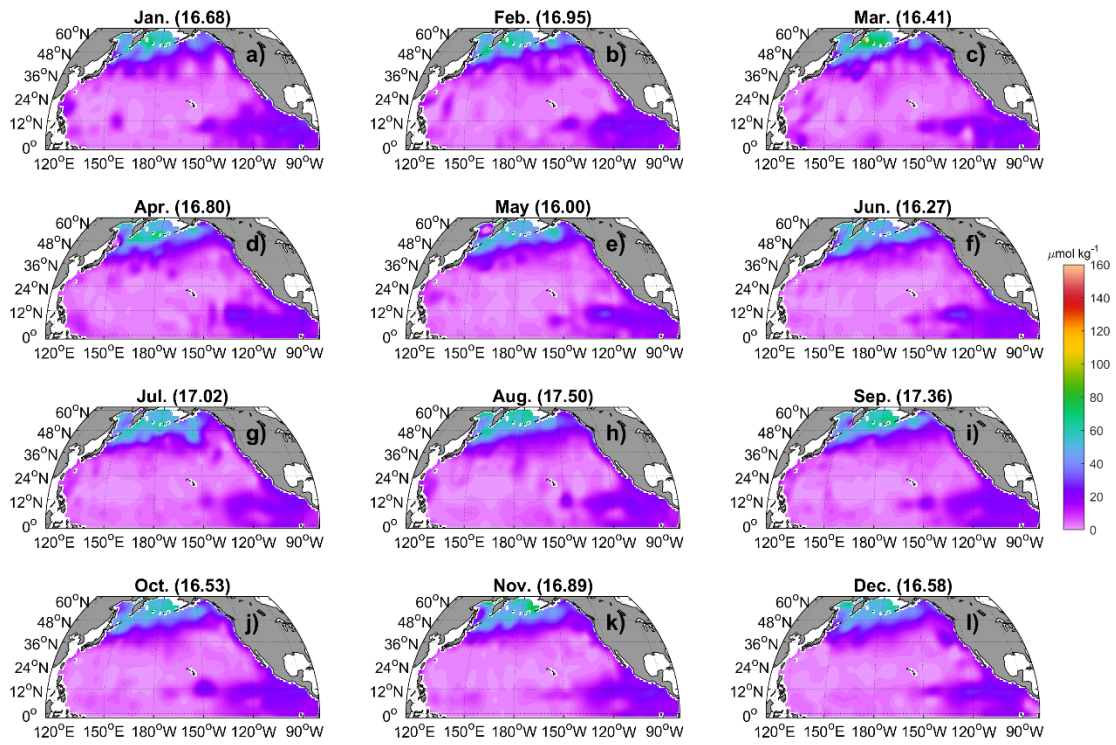


218

219 **Figure S42.** Similar to Fig. S36, but for Si(OH)_4 and at a depth of 5 m.

220

221

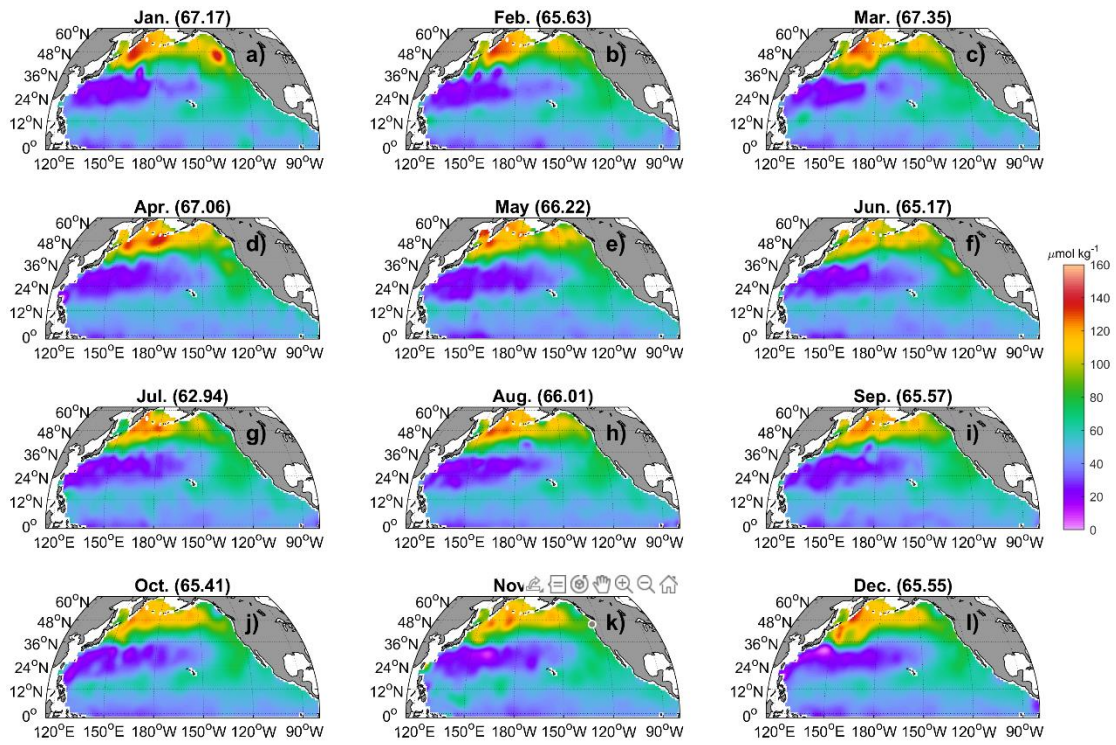


222

223 **Figure S43.** Similar to Fig. S36, but for Si(OH)_4 and at a depth of 100 m.

224

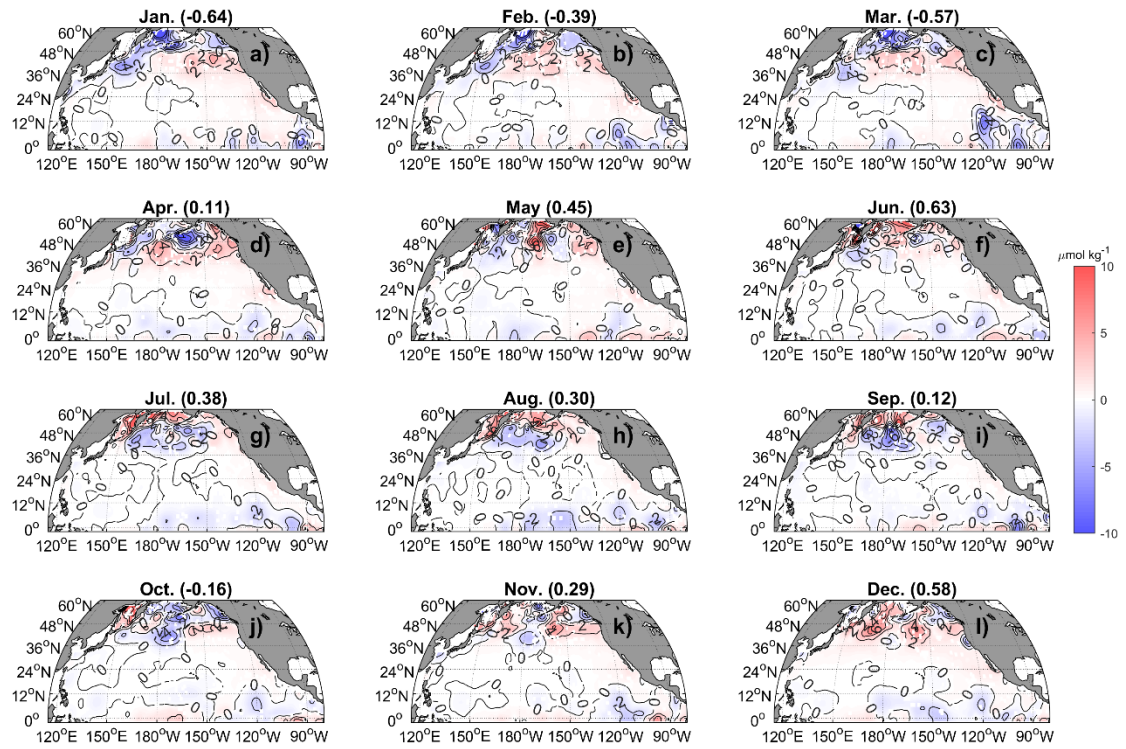
225



226

227 **Figure S44.** Similar to Fig. S36, but for Si(OH)_4 and at a depth of 500 m.

228



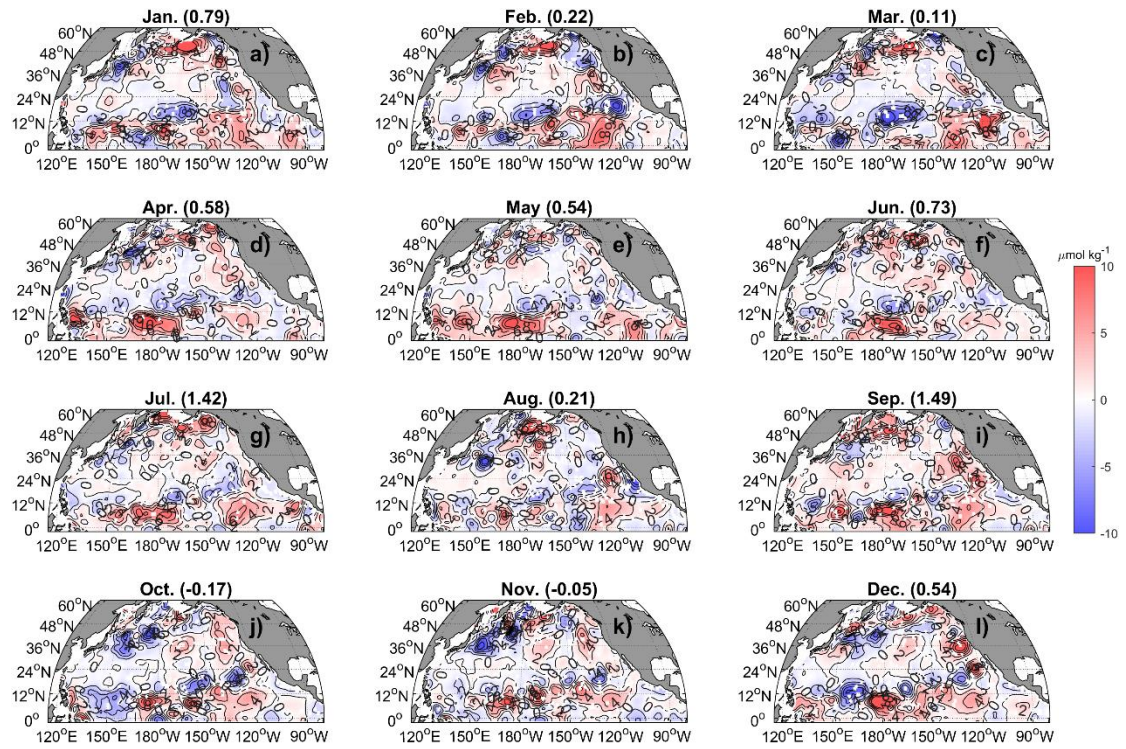
229

230 **Figure S45.** Difference between the predicted (Random Forest, RF) and World Ocean

231 Atlas (WOA) climatologies of NO_x^- ($\text{NO}_3^- + \text{NO}_2^-$) at 5 m depth (RF minus WOA).

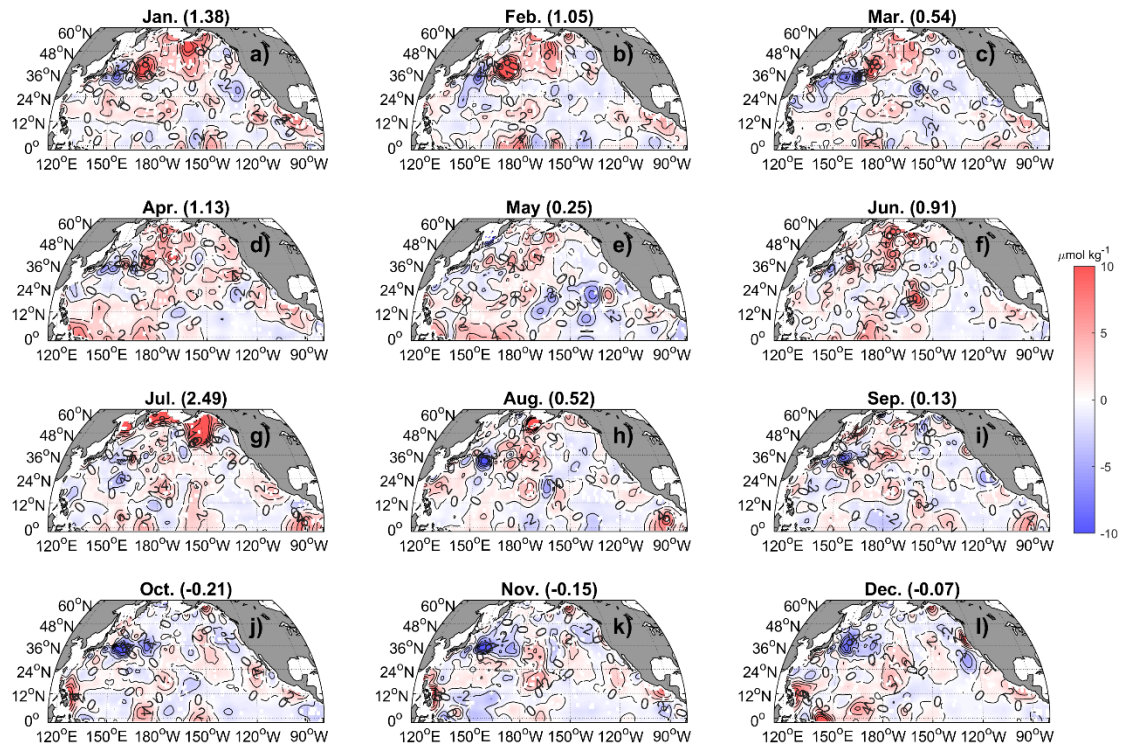
232

233



234
 235
 236
 237
 238

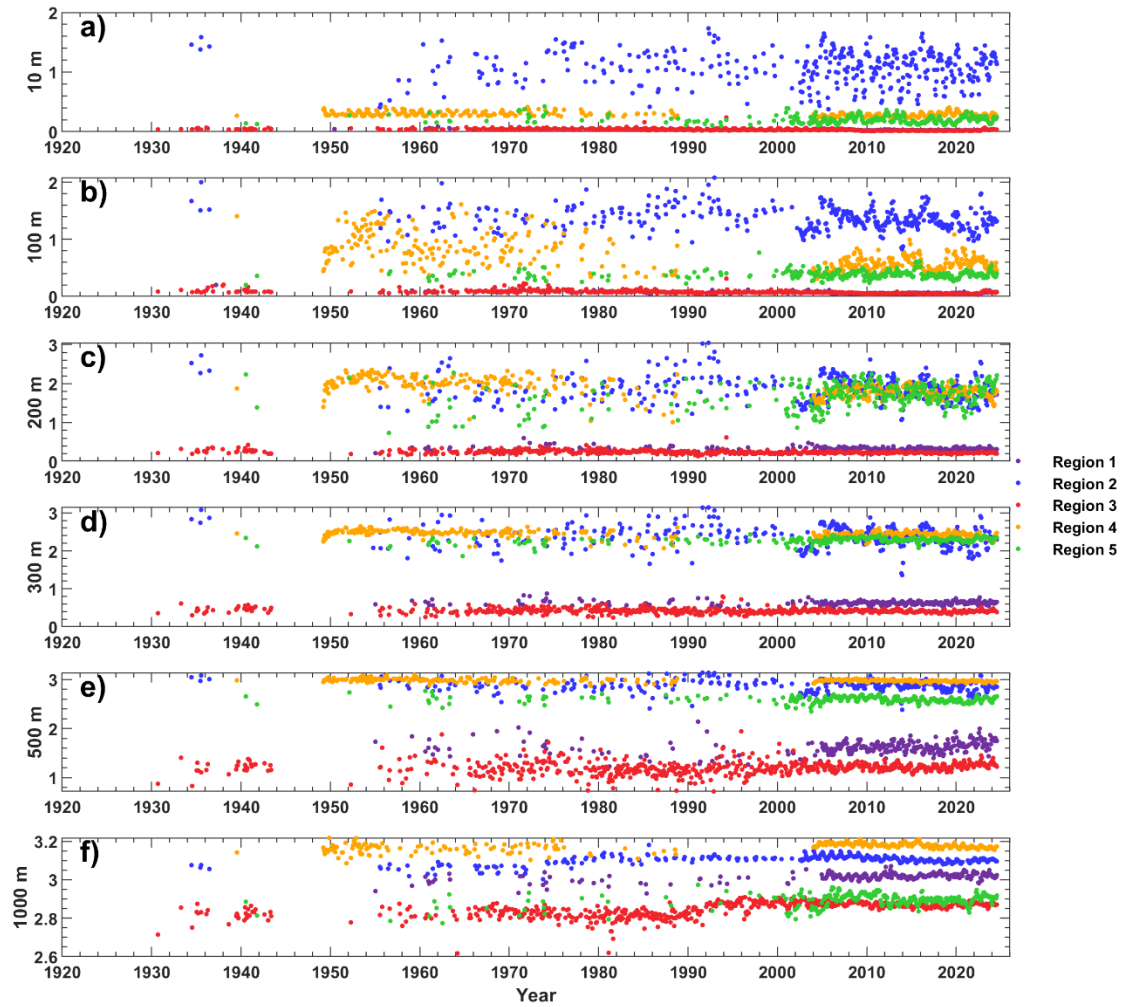
Figure S46. Same as Fig. S27, but for 100 m.



239

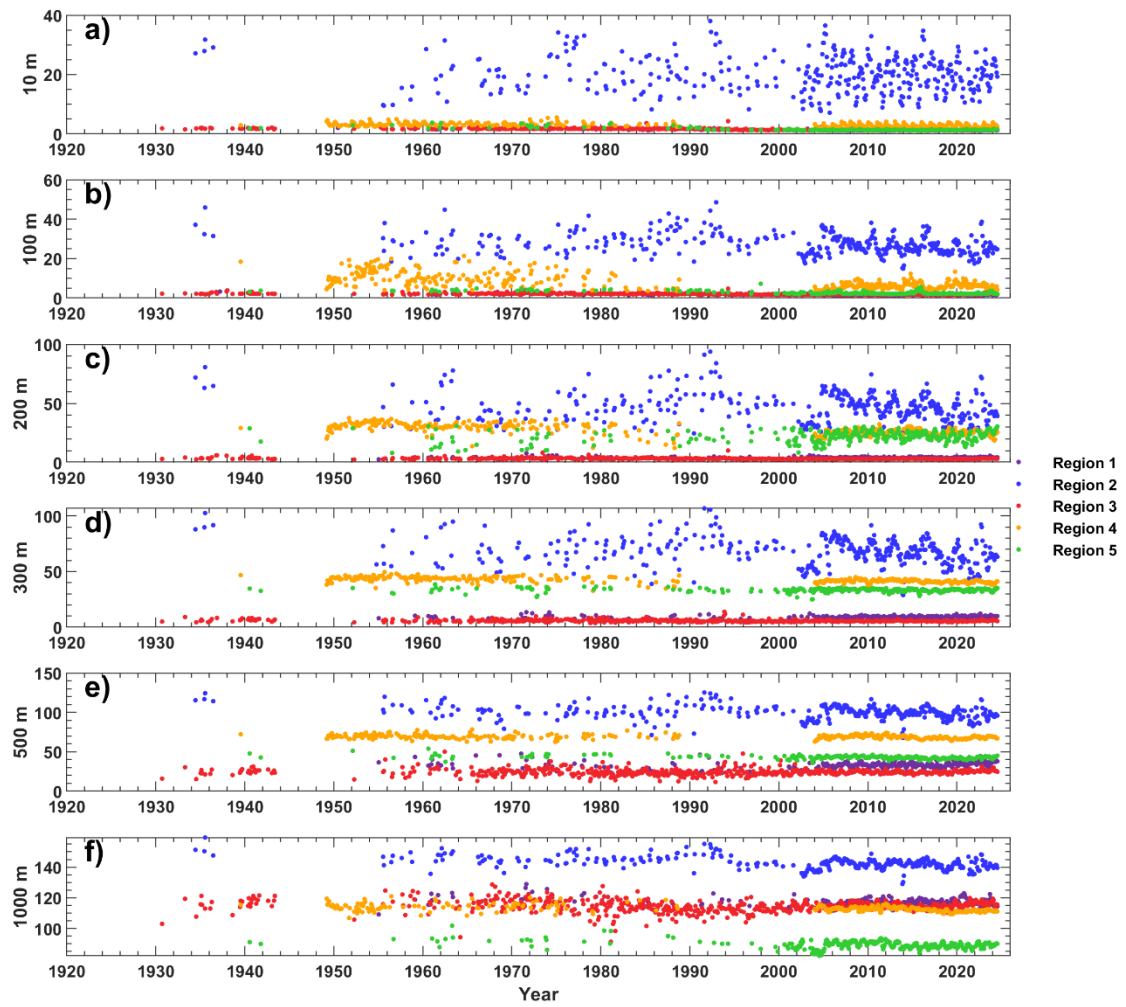
240

Figure S47. Same as Fig. S27, but for 500 m.



241

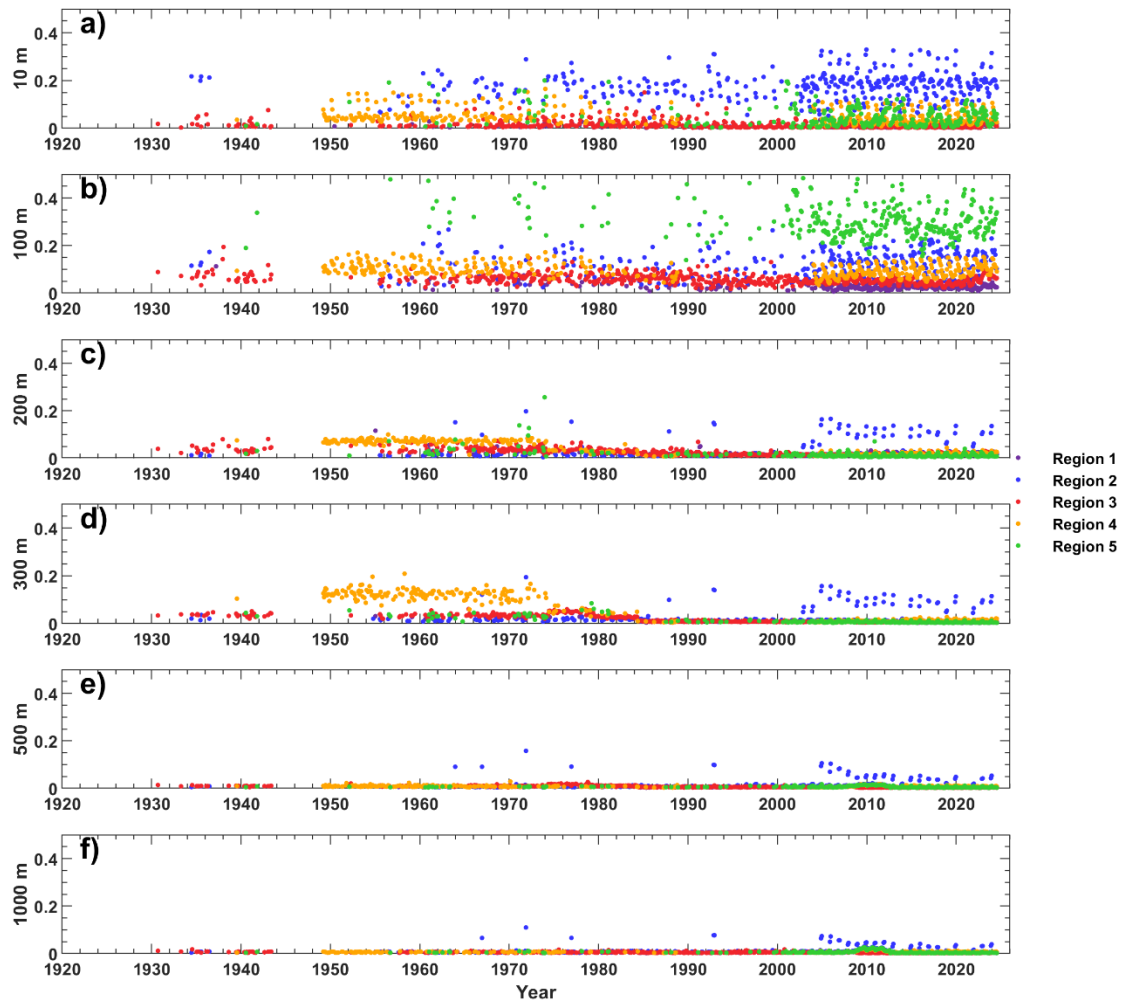
242 **Figure S48.** Time series of reconstructed DIP concentrations at 10 m (a), 100 m (b),
 243 200 m (c), 300 m (d), 500 m (e), and 1000 m (f) for regions 1 - 5 (see Fig. 16). Data
 244 were first binned by depth and region, then averaged by month.



245

246 **Figure S49.** Same as Figure S48, but for Si(OH)_4 .

247



248

249 **Figure S50.** Same as Figure S48, but for NO_2^- .

250

251