Earth System
Science
Data

Open Access

*Supplement of*

# A global monthly 3D field of seawater pH over 3 decades: a machine learning approach

**Guorong Zhong et al.**

*Correspondence to:* Xuegang Li (lixuegang@qdio.ac.cn) and Jinming Song (jmsong@qdio.ac.cn)

1    S1. Uncertainty and construction method of selected ocean products

2        A group of products related to the physical, chemical, and biological activities that

3    influence the ocean carbonate system were collected as potential pH predictors (Table

4    1). These products were constructed using different methods in previous research. The

5    seawater temperature and salinity product were constructed based on measurements

6    from the World Ocean Database (WOD) using the ensemble optimal interpolation

7    method with the dynamic ensemble (EnOI-DE) provided by CMIP5 historical

8    simulations (Cheng et al., 2016; Cheng et al., 2020). The temperature product was

9    claimed with an uncertainty of about ±0.05°C in the recent few decades, and the

10   uncertainty of salinity product was about ±0.001~±0.005 at different depths (present as

11   figures    in    Cheng    et    al.,    2016    and    Cheng    et    al.,    2020;

12   *https://journals.ametsoc.org/view/journals/clim/33/23/full-jcliD200366-f5.jpg*    and

13   *https://journals.ametsoc.org/view/journals/clim/29/15/full-jcli-d-15-0730.1-f8.jpg*).

14   The climatological Alk product was constructed from Global Ocean Data Analysis

15   Project version 2.2019 (GLODAPv2019) measurements using a neural network

16   (NNGv2) method, with the RMSE of 3–6.2 µmol kg$^{-1}$ (Broullón et al., 2019). The

17   climatological DIC product was constructed from GLODAPv2019 and the Lamont–

18   Doherty Earth Observatory (LDEO) datasets using a feedforward neural network

19   (dubbed NNGv2LDEO) method, with a RMSE of 3.6–13.2 µmol kg$^{-1}$ (Broullón et al.,

20   2020). The climatological dissolved oxygen, nitrate, phosphate, and silicate product

21   was constructed based on measurements from the World Ocean Database, using an

22   objective analysis method that generated a first-guess field and then carried out a

23   correction at all gridpoints as a distance-weighted mean of all gridpoint difference

24   values that lie within the area around the gridpoint defined by the influence radius

25   (Garcia et al., 2019a; Garcia et al., 2019b). The producer claimed an average DO bias

26   of 0.4±4.7 µmol kg$^{-1}$ below 500 m depth and 1.4±10.9 µmol kg$^{-1}$ above 500 m depth.

27   The average biases of nutrient concentration were -0.02±0.07 µmol kg$^{-1}$ for phosphate,

28   -0.22±0.95 µmol kg$^{-1}$ for nitrate, and -0.3±3.8 µmol kg$^{-1}$ for silicate below 500 m depth,

29   and were 0.01±0.12 µmol kg$^{-1}$ for phosphate, 0.2±1.8 µmol kg$^{-1}$ for nitrate, and 0.8±3.6

30   µmol kg$^{-1}$ for silicate above 500 m depth. The Sea surface height (SSH), mixed layer

31   depth (MLD), and W velocity of ocean current from the ECCO2 cube92 product were

32   constructed by least squares fit of a global full-depth-ocean and sea-ice configuration

33   of the Massachusetts Institute of Technology general circulation model to the available

34   satellite and in-situ data (Menemenlis et al., 2008). The basin-wide median bias error

35  of the MLD product is -6.6 m and the RMSE is 40 m, and the RMSE of the SSH product

36  is 9.2 cm. The ERA5 sea level pressure and surface pressure were constructed by the

37  Integrated Forecasting System (IFS) Cy41r2 model (Hersbach et al., 2020). The

38  standard deviation of ERA5 sea level pressure and surface pressure are within 1 hPa

39  and 0.8 hPa in the recent three decades. The NOAA Greenhouse Gas Marine Boundary

40  Layer Reference $xCO_2$ product is constructed by extending measurements from a subset

41  of sites from the NOAA Cooperative Global Air Sampling Network, with an uncertainty

42  within 1 μmol mol$^{-1}$ in most regions (Lan et al., 2023,

43  https://gml.noaa.gov/ccgg/mbl/mbl.html). The bi-monthly Multivariate El

44  Niño/Southern Oscillation index (MEI) was calculated by the first seasonally varying

45  principal component of six atmosphere–ocean (COADS) variable fields in the tropical

46  Pacific basin (Wolter et al., 2011). The Arctic Oscillation index was calculated as the

47  first leading mode from the Emperical Orthogonal Function analysis of monthly mean

48  height anomalies at 1000-hPa of the Northern Hemisphere or 700-hPa of the Southern

49  Hemisphere (CPC, 2002). The Southern Oscillation Index was calculated based on the

50  differences in air pressure anomaly between Tahiti and Darwin, Australia (CPC, 2005).

51  The specific uncertainty of these index products is not provided. The GEBCO global

52  bathymetric data was constructed using predicted depths based on the V32 gravity

53  model (Sandwell et al., 2019). The monthly surface ocean $p$CO$_2$ was constructed using

54  the SOM-FFNN method based on regional-specific predictors selected by the stepwise

55  FFNN algorithm, with a global RMSE of 17.99 μatm (Zhong et al., 2022). A

56  climatological $p$CO$_2$ product constructed by another SOM-FFNN model was also used,

57  with the RMSE of 18.3 μatm (Landschützer et al., 2020). The Euphotic Depth product

58  was constructed from remote sensing reflectance (RRS) data derived inherent optical

59  properties using Lee algorithm (Lee et al., 2007), with an average percentage error of

60  13.7%. The chlorophyll concentration product was constructed based on RRS at 2-4

61  wavelengths between 440 and 670 nm with an uncertainty of 1-2%, using the algorithm

62  of Hu et al. (2019) that combines an empirical band difference approach at low

63  chlorophyll concentrations with a band ratio approach at higher chlorophyll

64  concentrations. The photosynthetically available radiation (PAR) product was based on

65  the observed Top-of-Atmosphere (TOA) radiances in the 400-700nm range that do not

66  saturate over clouds using the algorithm of Frouin et al. (2002), with an RMSE of 3.6

67  Einstein/m$^2$/day. The product of the diffuse attenuation coefficient at 490 nm (Kd490)

68  was calculated using an empirical relationship derived from in situ measurements

69    of Kd490 and blue-to-green band ratios of RRS. The remote sensing reflectance

70    product was derived from ocean color sensors based on the spectral distribution of

71    reflected visible solar radiation upwelling from below the ocean surface and passing

72    through the sea-air interface. The total absorption and backscattering products were

73    calculated using the default global configuration of the Generalized Inherent Optical

74    Property (GIOP) model (Werdell et al., 2013).

75    S2. Validation of cross-boundary method

76        The cross-boundary method reduced the pH predicting error slightly, but improved

77    the discontinuity problem in the SOM boundary effectively (Figure S1 a-d). However,

78    the discontinuity problem was not completely solved and some boundary line existed

79    in the spatial distribution, especially in the deeper ocean that pH measurements are

80    much sparser (Figure S1 e-f). Even so, the performance of FFNN predicting was better

81    when the cross-boundary method was applied. Compared with taking average in the

82    boundary area, the cross-boundary method avoided subjectively modifying the

83    boundary data. Correspondingly, this method may not solve the discontinuity problem

84    perfectively in some situations. The cross-boundary method also decreased the

85    predicting error slightly in vertical boundary areas (2 layers near the mixed layer depth).

86    However, the improvement was minor in the vertical distribution, due to the natural

87    existing substantial vertical gradient of seawater pH near the mixed layer depth (Figure

88    S2). Overall, the cross-boundary method increases information about seawater pH

89    variation out of boundaries in the neural network learning process, reducing the outliers

90    near the SOM boundary and vertical boundary.

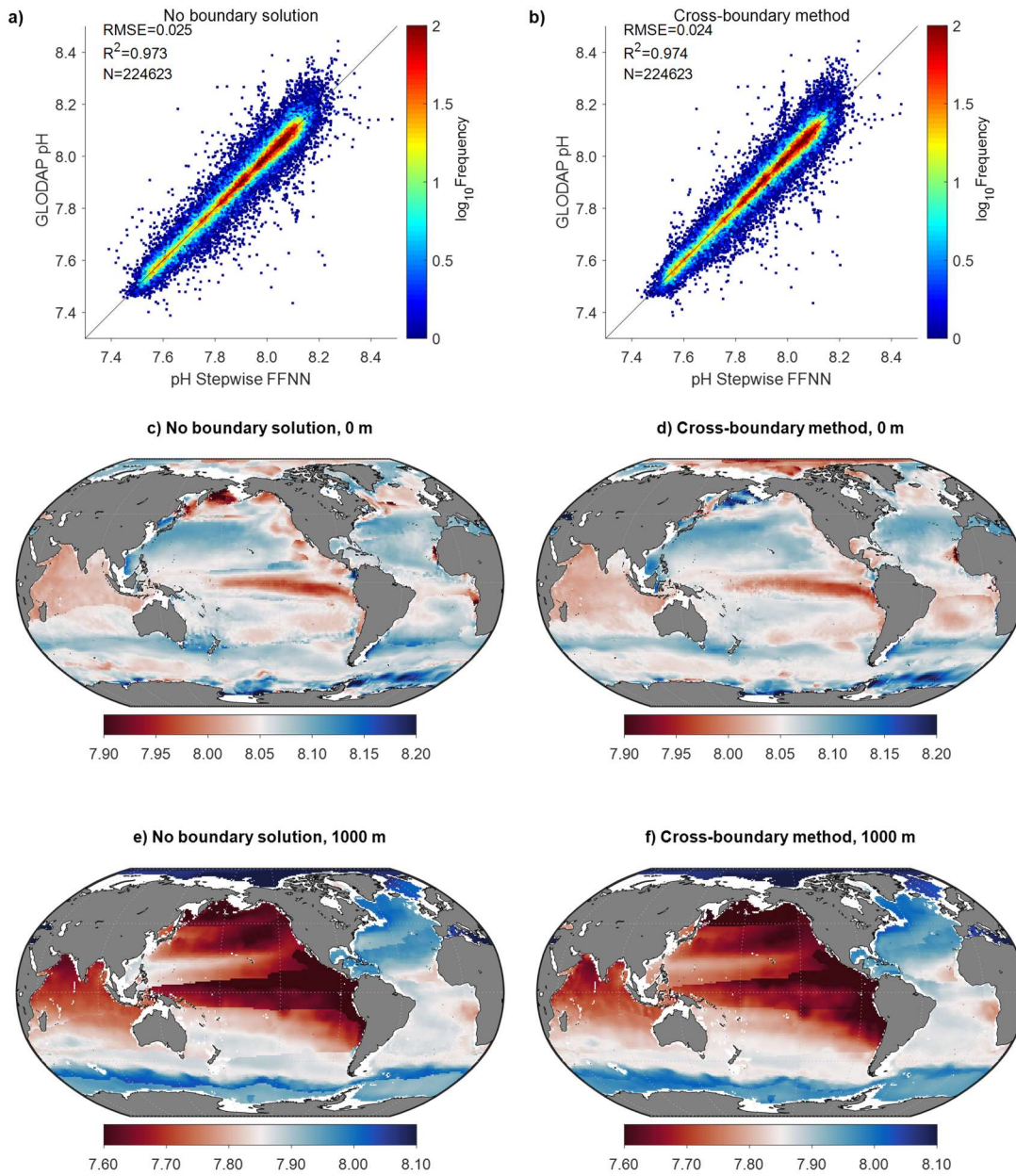91    S3. Comparison of performance between FFNNs training based on pH and $[H^+]$

92        Due to the logarithmic relationship between pH value and $[H^+]$ concentration,

93    results obtained from training FFNN with pH and from training FFNN with $[H^+]$ then

94    converting outputs into pH may differ. A comparison of predicting errors was conducted

95    between these two training methods. The results show a nearly consistent pH RMSE

96    between the FFNN training with pH and with $[H^+]$ (Figure 7). As the pH measurements

97    of all GLODAP samples are closer to a normal distribution than the $[H^+]$, the predicting

98    error was slightly lower in most regions when the FFNN was trained with pH, but the

99    difference in predicting errors was extremely small. In addition, the FFNN trained using

100   $[H^+]$ occasionally produced negative $[H^+]$ in regions with extremely low $[H^+]$.

101   Therefore, it is better to train FFNN using pH rather than using $[H^+]$ in the

102   reconstruction process of the pH product.

103　　　　The distribution patterns of regional pH RMSE and [H$^+$] RMSE are inconsistent

104　　whenever the FFNN was trained using pH or [H$^+$]. In fact, the pH RMSE of the

105　　intermediate layer in regions such as the subarctic North Pacific and the equatorial

106　　Pacific is significantly lower than that in the intermediate layer of the Arctic Ocean, but

107　　their [H$^+$] RMSE is higher than that of the intermediate layer in the Arctic Ocean (Figure

108　　7a and 7b). This is caused by the effect of the logarithmic relationship. If the pH values

109　　are different for the same pH RMSE, the corresponding [H$^+$] RMSE will be different.

110　　Therefore, the uncertainty of the pH product is calculated based on the [H$^+$] RMSE and

111　　pH value, rather than solely based on the pH RMSE.

112

113

119
120

121   **Figure S2. Validation of cross-boundary method for pH predicting in the SOM boundary**. a-
122   b): comparison of FFNN predicted pH with GLODAP in all SOM boundary areas; c-f):
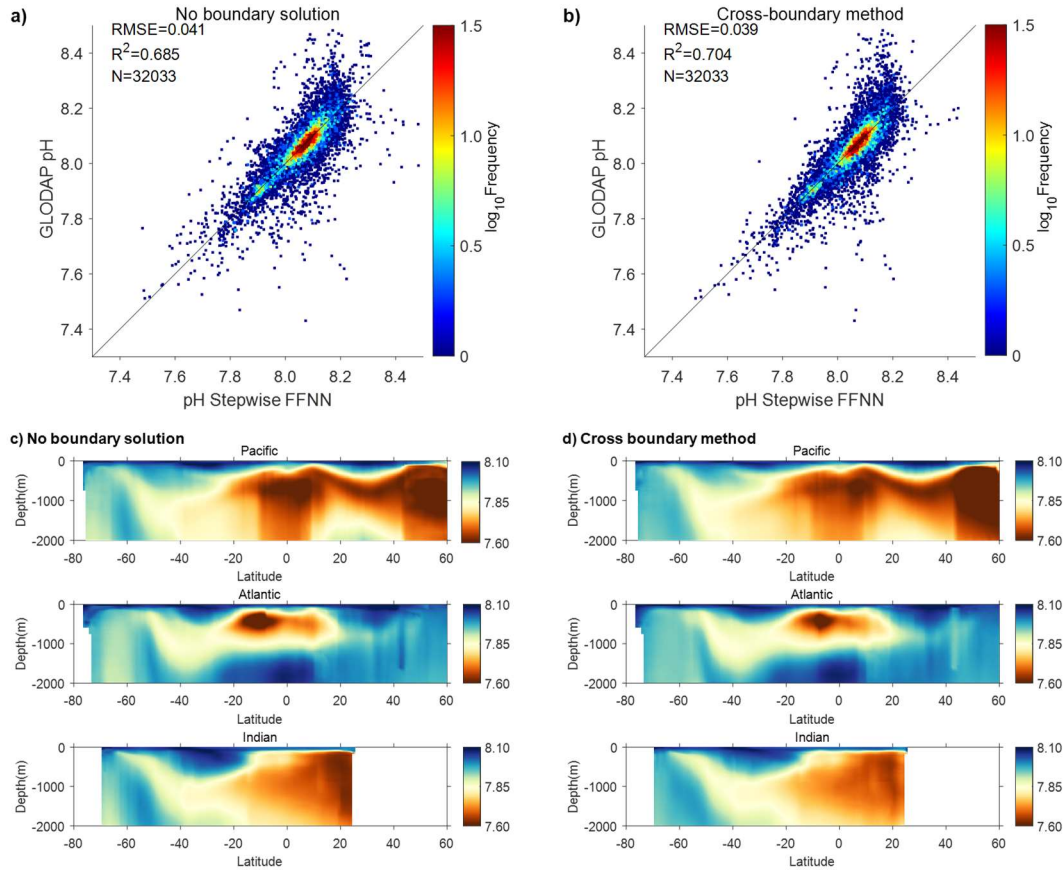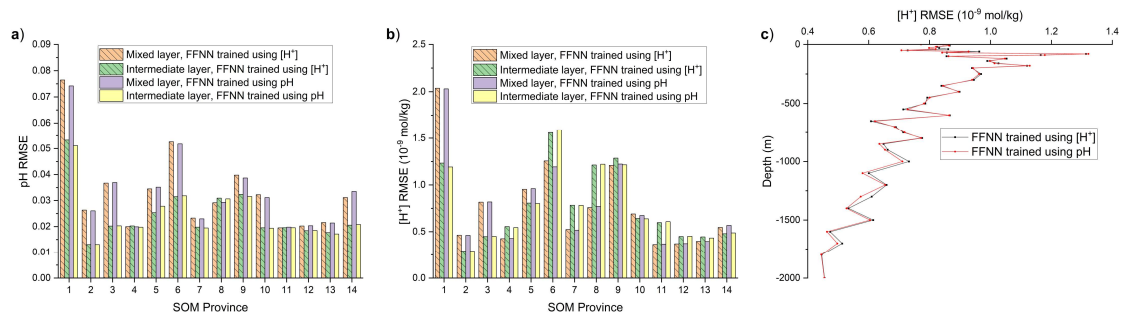123   comparison of spatial distribution at 0 m and 1000 m in January 2020.



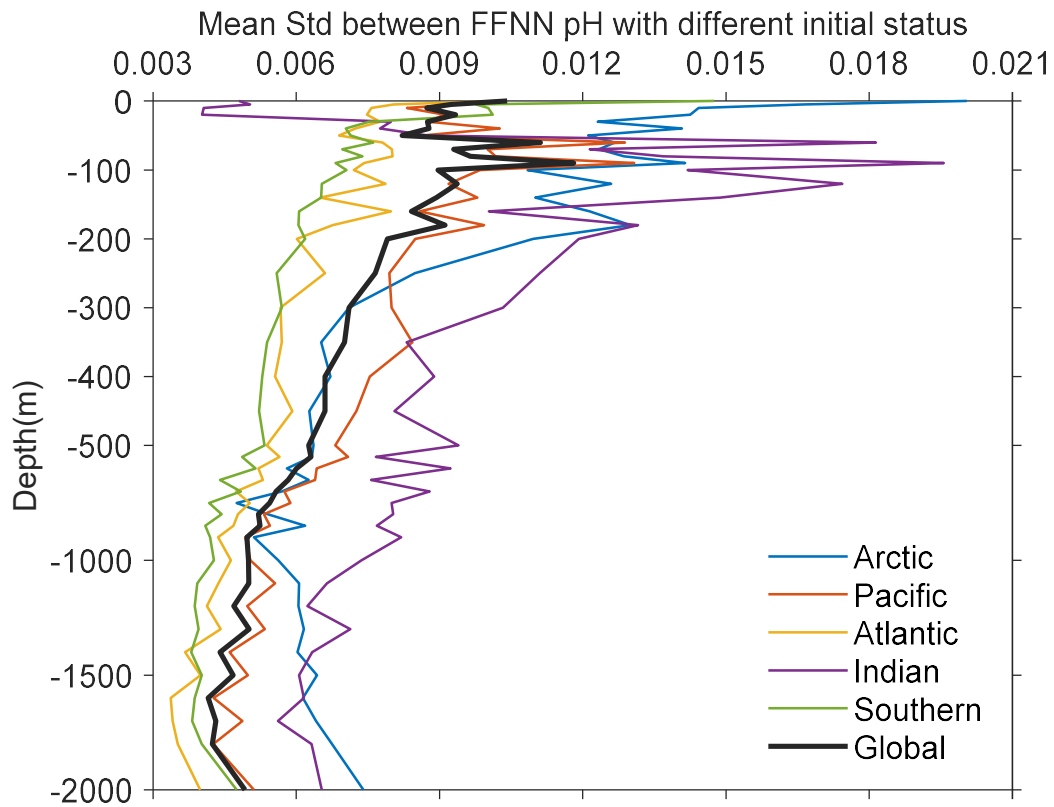124

125

126
127

132

141

142

143

144 **Figure S5. Mean standard deviation between FFNN pH with different initial status.**


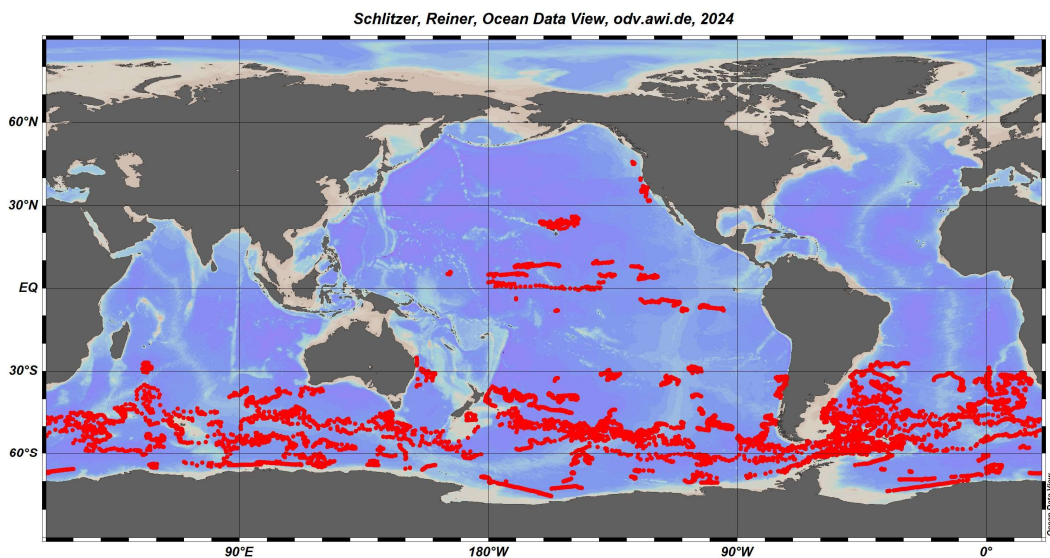
145

146

147

148 **Figure S6. Station map of used delayed-mode BGC-Argo pH-adjusted data with quality**
149 **control flag 1.**



150

151

152

153

154

**Table S1. Predictors selected by the stepwise FFNN algorithm in the Mixed layer for period before August 2002.** The predictors are arranged in order of relative importance, with the variables listed at the front of each province being more effective in reducing predicting errors when used as pH predictors.

| Province | FFNN neurons | Predictors |
|---|---|---|
| P5 Equatorial Atlantic | 25 | Phosphate, Temp, SLP, DIC, $P_{surf}$, TA, $pCO_2$, $W_{vel}$(in-situ), DO |
| P8 Equatorial Pacific | 10 | $pCO_2$, Depth, sLat, Temp, Sal, DIC, $W_{vel}$(in-situ), Nitrate |
| P10 Subtropical South Atlantic | 20 | $pCO_2$, Silicate, Nitrate, $W_{vel}$(65m), $W_{vel}$(in-situ), $W_{vel}$(195m) |
| P11 Subtropical South Pacific | 10 | Phosphate, $pCO_2$, Depth, sLat, Silicate, $pCO_{2\ clim}$, $W_{vel}$(5m), $W_{vel}$(105m) |

**References mentioned in supplementary text:**

Broullón, D., Pérez, F. F., Velo, A., Hoppema, M., Olsen, A., Takahashi, T., Key, R. M., Tanhua, T., González-Dávila, M., Jeansson, E., Kozyr, A., and van Heuven, S. M. A. C.: A global monthly climatology of total alkalinity: a neural network approach, Earth Syst. Sci. Data, 11, 1109–1127, https://doi.org/10.5194/essd-11-1109-2019, 2019.

Broullón, D., Pérez, F. F., Velo, A., Hoppema, M., Olsen, A., Takahashi, T., Key, R. M., Tanhua, T., Santana-Casiano, J. M., and Kozyr, A.: A global monthly climatology of oceanic total dissolved inorganic carbon: a neural network approach, Earth Syst. Sci. Data, 12, 1725–1743, https://doi.org/10.5194/essd-12-1725-2020, 2020.

Cheng, L. and Zhu, J.: Benefits of CMIP5 multimodel ensemble in reconstructing historical ocean subsurface temperature variations, J. Clim., 29, 5393-5416, https://doi.org/10.1175/JCLI-D-15-0730.1, 2016.

Cheng, L., Trenberth, K. E., Gruber, N., Abraham, J. P., Fasullo, J. T., Li, G., Mann, M. E., Zhao, X., and Zhu, J.: Improved estimates of changes in upper ocean salinity and the hydrological cycle, J. Clim., 33, 10357-10381, https://doi.org/10.1175/JCLI-D-20-0366.1, 2020.

Climate Prediction Center: Daily Arctic Oscillation Index [data set], https://www.cpc.ncep.noaa.gov/products/precip/CWlink/daily_ao_index/ao_index.html, 2002.

Climate Prediction Center: Southern Oscillation Index [data set], https://www.cpc.ncep.noaa.gov/products/analysis_monitoring/ensocycle/soi.shtml, 2005.

Frouin, R., Franz, B. A., and Werdell, P. J.: The SeaWiFS PAR product. ,In: S.B. Hooker and E.R. Firestone, Algorithm Updates for the Fourth SeaWiFS Data

Reprocessing, NASA Tech. Memo, 2003-206892, Volume 22, NASA Goddard Space Flight Center, Greenbelt, Maryland, 46-50, 2002.

Garcia, H. E., Weathers, K. W., Paver, C. R., Smolyar, I., Boyer, T. P., Locarnini, R. A., Zweng, M. M., Mishonov, A. V., Baranova, O. K., Seidov, D., and Reagan, J. R.: World Ocean Atlas 2018, Volume 3: Dissolved Oxygen, Apparent Oxygen Utilization, and Dissolved Oxygen Saturation, edited by: Mishonov, A., NOAA Atlas NESDIS 83, 38 pp., https://www.nodc.noaa.gov/OC5/woa18/pubwoa18.htm, 2019a.

Garcia, H. E., Weathers, K. W., Paver, C. R., Smolyar, I., Boyer, T. P., Locarnini, R. A., Zweng, M. M., Mishonov, A. V., Baranova, O. K., Seidov, D., and Reagan, J. R.: World Ocean Atlas 2018. Vol. 4: Dissolved Inorganic Nutrients (phosphate, nitrate and nitrate+nitrite, silicate). A. Mishonov Technical Editor, NOAA Atlas NESDIS 84, 35 pp., https://archimer.ifremer.fr/doc/00651/76336/, 2019b.

Hersbach, H., Bell, B., Berrisford, P., Hirahara, S., Horányi, A., Muñoz-Sabater, J., Nicolas J., Peubey C., Radu R., Schepers D., Simmons A., Soci C., Abdalla S., Abellan X., Balsamo G., Bechtold P., Biavati G., Bidlot J., Bonavita M., Chiara G. D., Dahlgren P., Dee D., Diamantakis M., Dragani R., Flemming J., Forbes R., Fuentes M., Geer A., Haimberger L., Healy S., Hogan R. J., Hólm E., Janisková M., Keeley S., Laloyaux P., Lopez P., Lupu C., Radnoti G., Rosnay P. D., Rozum I., Vamborg F., Villaume S., and Thépaut, J. N.: The ERA5 global reanalysis, Q. J. R. Meteorol. Soc., 146, 1999-2049, https://doi.org/10.1002/qj.3803, 2020.

Hu, C., Feng, L., Lee, Z., Franz, B. A., Bailey, S. W., Werdell, P. J., and Proctor, C. W.: Improving satellite global chlorophyll a data products through algorithm refinement and data recovery, J. Geophys. Res.-Oceans, 124(3), 1524-1543, https://doi.org/10.1029/2019JC014941, 2019.

Lan, X., Tans, P., Thoning, K., and NOAA Global Monitoring Laboratory: NOAA Greenhouse Gas Marine Boundary Layer Reference - $CO_2$, NOAA GML [Data set], https://doi.org/10.15138/DVNP-F961, 2023.

Landschützer, P., Laruelle, G. G., Roobaert, A., and Regnier, P.: A uniform $p\mathrm{CO_2}$ climatology combining open and coastal oceans, Earth Syst. Sci. Data, 12, 2537–2553, https://doi.org/10.5194/essd-12-2537-2020, 2020.

Lee, Z., Weidemann, A., Kindle, J., Arnone, R., Carder, K. L., and Davis, C.: Euphotic zone depth: Its derivation and implication to ocean-color remote sensing, J. Geophys. Res., 112, C3, https://doi.org/10.1029/2006JC003802, 2007.

Menemenlis, D., Campin, J. M., Heimbach, P., Hill, C., Lee, T., Nguyen, A., Schodlok, M., and Zhang, H.: ECCO2: High resolution global ocean and sea ice data synthesis, *Mercat. Ocean Q. Newsl,* 31, 13-21, 2008.

Sandwell, D. T., Harper, H., Tozer, B., and Smith, W. H.: Gravity field recovery from geodetic altimeter missions, Adv. Space Res., 68(2), 1059-1072, https://doi.org/10.1016/j.asr.2019.09.011, 2021.

230 Werdell, P. J., Franz, B. A., Bailey, S. W., Feldman, G. C., Boss, E., Brando, V. E.,
231     Dowell M., Hirata T., Lavender S. J., Lee, Z., Loisel H., Maritorena S., Mélin F.,
232     Moore T. S., Smyth T. J., Antoine D., Devred E., d'Andon O. H. F., and Mangin, A.:
233     Generalized ocean color inversion model for retrieving marine inherent optical
234     properties. Appl. Optics, 52(10), 2019-2037, https://doi.org/10.1364/AO.52.002019,
235     2013.
236 Wolter, K. and Timlin, M. S.: El Niño/Southern Oscillation behaviour since 1871 as
237     diagnosed in an extended multivariate ENSO index (MEI. ext), Int. J. Climatol., 31,
238     1074-1087, https://doi.org/10.1002/joc.2336, 2011.
239 Zhong, G., Li, X., Song, J., Qu, B., Wang, F., Wang, Y., Zhang, B., Sun, X., Zhang,
240     W., Wang, Z., Ma, J., Yuan, H., and Duan, L.: Reconstruction of global surface
241     ocean $p$CO$_2$ using region-specific predictors based on a stepwise FFNN regression
242     algorithm, Biogeosciences, 19, 845–859, https://doi.org/10.5194/bg-19-845-2022,
243     2022.
244