



Retrieving monthly and interannual pH_T on the East China Sea shelf using an artificial neural network: ANN- pH_T -v1

Xiaoshuang Li^{1,2}, Richard Bellerby^{1,2}, Jianzhong Ge¹, Philip Wallhead², Jing Liu¹, and Anqiang Yang¹

¹State Key Laboratory of Estuarine and Coastal Research, East China Normal University, Shanghai, 200241, China

5 ²Norwegian Institute for Water Research, Bergen, 5006, Norway

Correspondence to: Richard Bellerby (Richard.Bellerby@niva.no)

Abstract. While our understanding of pH dynamics has strongly progressed for open ocean regions, for marginal seas such as the East China Sea (ECS) progress has been constrained by limited observations and complex interactions between biological, physical, and chemical processes. Seawater pH is a very valuable oceanographic variable but not always measured using high quality instrumentation and according to standard practices. In order to predict water column total scale pH (pH_T) and enhance our understanding of the seasonal variability of pH_T on the ECS shelf, an artificial neural network (ANN) model was developed using 11 cruise datasets from 2013 to 2017 with coincident observations of pH_T , temperature (T), salinity (S), dissolved oxygen (DO), nitrate (N), phosphate (P) and silicate (Si) together with sampling position and time. The reliability of the ANN model was evaluated using independent observations from 3 cruises in 2018, and showed a root mean square error accuracy of 0.04. A weight analysis of the ANN model variables suggested that DO, S, T were the most important predictor variables. Monthly water column pH_T for the period 2000-2016 was retrieved using T, S, DO, N, P, and Si from the Changjiang Biology Finite-Volume Coastal Ocean Model (FVCOM).

1 Introduction

Atmospheric carbon dioxide (CO_2) levels have increased by nearly 46%, from approximately 278 ppm (parts per million) in 1750 (Ciais et al., 2013) to 405 ppm in 2017 (Le Quéré et al., 2018). The oceans have absorbed approximately 48% of the anthropogenic CO_2 emissions (Sabine et al., 2004), resulting in decreasing long-term pH trends of ~ 0.02 decade⁻¹ in open ocean waters (e.g., Dore et al., 2009; González-Dávila et al., 2010; Bates et al., 2014; Lauvset et al., 2015). While a gradual decrease in pH is a predictable open ocean response to elevated anthropogenic CO_2 emissions, the seasonal changes and long-term trends in pH in coastal seas have not been fully understood due to the lack of long-term pH data and complexity of coastal systems. In this context, the development of approaches to predict carbonate chemistry parameters in coastal regions may assist both the management of local water quality and our wider understanding of the ocean carbon cycle.

Many attempts have been made to predict seawater pH by developing empirical relationships between pH and environmental variables, such as temperature (T) (Juraneck et al., 2011), salinity (S) (Williams et al., 2016), dissolved oxygen (DO) (e.g., Juraneck et al., 2011; Sauzède et al., 2017), nutrients (e.g., Williams et al., 2016; Carter et al., 2016, 2018), and longitude, latitude (Sauzède et al., 2017). Compared with traditional empirical methods, artificial neural networks (ANNs) have shown improved performance (Chen et al., 2017). ANNs have been proposed as powerful tools for modelling uncertain and complex systems such as ecosystems and environmental assessment (e.g., Olden and Jackson, 2002; Olden et al., 2004; Uusitalo, 2007; Raitso et al., 2008). Their main advantage compared with multiple linear regression (MLR) models is that they do not require an a priori model but rather “learn” the model from existing data. ANNs have been used for the retrieval of the partial pressure of carbon dioxide (pCO_2) (e.g., Friedrich and Oschlies, 2009; Laruelle et al., 2017), total alkalinity (e.g., Velo et al., 2013; Bostock et al., 2013; Sasse et al., 2013), total dissolved inorganic carbon (e.g., Bostock et al., 2013; Sasse et al., 2013), and phytoplankton functional types (e.g., Raitso et al., 2008; Palacz et al., 2013). However, these studies mainly focus on the open ocean; relatively few studies have focused on coastal seas, perhaps because of the complexity and heterogeneity of the



40 continental shelves. Alin et al. (2012) developed an MLR model to reconstruct pH in the southern California Current System, while Moore-Maley et al. (2016) evaluated the interannual variability of near-surface pH using a one dimensional, biophysical, mixing layer model in the Strait of Georgia. To our knowledge, no empirical relationship for pH has yet been established for the ECS.

45 The ECS is the largest marginal sea in the western North Pacific Ocean and receives massive terrestrial inputs from the Changjiang River (Gong et al., 1996). The spatial and temporal distributions of the carbonate system have been investigated in the ECS (e.g., Chou et al., 2009; Cao et al., 2011; Qu et al., 2015), and were found to largely reflect the distributions of various water masses. The pattern of carbon sources and sinks exhibits substantial seasonal variation (Guo et al., 2015), and the ECS is generally considered as a sink of atmospheric CO₂ throughout the year except in fall (e.g., Shim et al., 2007; Zhai
50 and Dai, 2009). A mechanistic semi-analytical algorithm (MeSAA) was developed to study pCO₂ variations in response to various controlling mechanisms during summertime (Bai et al., 2015). However, the seasonal variability of pH has been very little studied in the ECS, mainly due to the limited observational coverage and irregular variability caused by seasonal fluctuations of the Changjiang River discharge and anthropogenic processes. Developing methods to extend the seasonal coverage of pH data may thus help to improve our understanding of the ocean carbon cycle in the ECS.

55

This paper is structured as follows: section 2 describes the cruise data used to build the ANN model; section 3 shows the ANN model performance and predictor variable importance, as well as an application to retrieve monthly pH for the period 2000-2016 on the ECS shelf using the monthly temperature, salinity, dissolved oxygen, nitrate, phosphate and silicate from the Changjiang Biology Finite-Volume Coastal Ocean Model (FVCOM). Conclusions and perspectives are summarized in the last
60 section.

2 Data and method

2.1 Data

Eleven cruises were conducted on the ECS shelf from 2013 to 2017 (Fig. 1). Ten of these were carried out during the “Fund Committee Innovation Group Project” (Y22323101B), the summer cruise from 17 to 28 August 2013, 10 to 17 July 2014, 9
65 to 20 July 2015, 4 to 28 July 2016, 20 to 30 July 2017, the winter cruise from 21 to 28 February 2014, 15 to 28 February 2017, the spring cruise from 4 to 20 March 2013, 11 to 21 March 2015, 7 to 19 March 2016; the remaining cruise was carried out on the ECS shelf during 12–24 May 2017. Water samples were collected at three or four different depths during all cruises.

T and S profiles were obtained directly using a conductivity temperature-depth/pressure (CTD) recorders (SBE 25plus or
70 911plus). Measurement of DO followed the Winkler procedure, as described previously by Zhai et al. (2014). Nutrients samples were first filtered with 0.45 μm Whatman GF/F membrane, then stored in 250 mL HDPE bottles until chemical analysis. Nitrate (N), phosphate (P) and silicate (Si) were determined using a segmented flow analyzer (Model: Skalar SAN^{PLUS}, Netherlands) with a precision < 3% (Zhang et al., 2007), the detection limits are 0.14 μM for N, 0.06 μM for P, and 0.07 μM for Si. During the May 2017 survey, pH samples were stored in 500 mL high-quality borosilicate glass bottles without filtering
75 and poisoned by addition of 200 μL saturated HgCl₂ solution until measured in the lab. The pH_T (total scale) was measured using an Automated Flow-through system for Embedded Spectrophotometry (AFtes) with a precision of 0.0005 pH_T unit and uncertainty of < 0.003 (Reggiani et al., 2016). During all other cruises, pH samples were stored in 140 mL brown borosilicate glass bottles and poisoned by addition of 50 μL saturated HgCl₂ solution. Three traceable pH buffers were used including NIST (National Institute of Standards and Technology) buffers pH = 4.00, 7.02, 10.09. As described by Zhai et al. (2012,
80 2014), we converted it into total scale and the overall accuracy of the pH_T data was estimated as 0.01. We omitted data points



where one or more other physical variables were missing. The final number of data used by the ANN model was 1854, and the distribution of the sampling sites from the 11 cruises is shown in Fig. 1. More detailed information on the field survey can be found in Table 1.

2.2 Artificial neural network development

85 The ANN we used is a feed-forward multilayer perceptron (Tamura and Tateishi, 1997) with two hidden layers. The neurons of each layer are connected with the neurons of the previous layer and the next layer by weights (Fig. 2). The coefficients of the weight matrix are iteratively tuned in the training step. Here we used the back-propagation conjugate-gradient technique (Hornik et al., 1989). In order to avoid overfitting, a ten-fold cross-validation was used to assess model prediction accuracy. Here, the cruise data were randomly divided into ten equal subsamples. One subsample was used as the independent validation
90 data (10% of all cruise data), which was always excluded from training, and the nine remaining subsamples were together used as training data (90% of all cruise data). The training data were further divided randomly into a training set (70% of training data), validation set (15% of training data), and testing set (15% of training data). The training set was used for computing the gradient and updating the network weights and biases. The validation set was used to monitor the error and control model stop during the training process. The testing set was used to monitor whether the model was over-matched (Palacz et al., 2013). We
95 compared performances in predicting the independent validation data from the ten-fold cross-validation and selected the optimal model based on the lowest root mean square error (RMSE). In our study, calculations were done in the MathWorks Matlab environment, using the Deep Learning Toolbox.

Different combinations of input variables were tested to choose the optimal architecture of the ANN model (Table 2). As
100 shown in Fig. 2, input variables include longitude, latitude, month, T, S, DO, N, P and Si. We selected these variables as principal inputs for the following reasons: the carbonate system thermodynamic relationships depend on T and S (Lueker et al., 2000); DO was expected to vary with pH and there was a tight positive link between DO and pH (Wootton et al., 2012) because of the role of photosynthesis and respiration in removing or generating CO₂ in the water; nutrients influence phytoplankton growth and abundance, which might increase organic carbon fixation, increasing inorganic carbon uptake and
105 increasing pH (Wootton et al., 2008, 2012). We found geographical information to be a powerful addition in improving the skill of the method (see Table 2), allowing the network to learn spatio-temporal patterns that could not be explained by other input variables (Sasse et al., 2013). The number of neurons in the two hidden layers was tested, varying between 1 and 100 for the first hidden layer and between 1 and 50 for the second hidden layer. The optimal architecture was composed of two hidden layers with 40 neurons in the first and 16 neurons in the second.

110

In order to avoid bias towards high-value inputs/outputs and to eliminate the dimensional influence of the data, all data used by the ANN model were normalized using the following equation (e.g., Sauzède et al., 2015, 2016):

$$x_{i,j} = \frac{2}{3} * \frac{x_{i,j} - \text{mean}(x_{i,j})}{\sigma(x_{i,j})} \quad (1)$$

with σ the standard deviation of the considered input variables or output variable pH_T. Similar to the approach of Sauzède et
115 al. (2015, 2016), the longitude and month input variables were transformed as follows to account for the periodicity:

$$\text{slongitude} = \sin\left(\frac{\text{Lon} + \pi}{180}\right), \quad \text{clongitude} = \cos\left(\frac{\text{Lon} + \pi}{180}\right) \quad (2)$$

$$\text{smmonth} = \sin\left(\frac{\text{month} + \pi}{6}\right), \quad \text{cmmonth} = \cos\left(\frac{\text{month} + \pi}{6}\right) \quad (3)$$

The latitude variable was transformed into the range of the sigmoid function by dividing by 90 (Sauzède et al., 2015), then normalized using (1).



120 **3 Result and discussion**

3.1 The ANN model performance

To evaluate the performance of the ANN model, we compared model simulated pH_T (pH_T^M) with corresponding observations (pH_T^O) using several statistical indices, including the mean absolute error (MAE), the coefficient of determination (R^2), and the root mean squared error (RMSE). The model simulated pH_T with $\text{RMSE} = 0.04$ and $R^2 = 0.88$ for the training data (90%
125 of all data, Fig 3a), and predicted pH_T with $\text{RMSE} = 0.03$ and $R^2 = 0.93$ for the independent validation data (10% of all data, Fig 3b). The distributions of the differences ($\text{pH}_T^M - \text{pH}_T^O$) were approximately normal with no obvious outliers (Fig. 4).

In order to further explore where the ANN model may lead to large errors, we plotted distributions of differences ($\text{pH}_T^M - \text{pH}_T^O$) with respect to the longitude and latitude (Fig. 5). The points with large errors are mainly concentrated in the longitude
130 range [122.5 °E, 123 °E] and the latitude range [32 °N, 32.5 °N], in an area strongly influenced by the Changjiang Dilute Water (CDW). The reduced performance of the ANN model here can be primarily attributed to seasonal oscillations of the Changjiang River discharge (Dai and Trenberth, 2002). Although the RMSE for pH_T we obtained here was higher than obtained in some previous studies (e.g., Juranek et al., 2011; Williams et al., 2016; Sauzède et al., 2017), their research regions were open ocean regions, not coastal seas. For example, Juranek et al. (2011) developed empirical algorithms to estimate pH with RMSE of
135 0.018 for data between 30-500 m in the NE subarctic Pacific; Williams et al. (2016) also developed empirical algorithms to predict pH with RMSE of 0.01 in the Southern Ocean; Sauzède et al. (2017) developed a neural network method to estimate pH with RMSE of 0.02 in the global ocean. However, coastal seas tend to show greater temporal and spatial variability than open oceans. Alin et al. (2012) developed a MLR approach to reconstruct pH with RMSE of 0.024 in the southern California Current System. Zhai et al. (2014) compared the field-measured pH with calculated pH from measured total alkalinity and
140 dissolved inorganic carbon using the program CO2SYS.xls (Pelletier et al., 2011) and obtained discrepancies with standard deviation 0.05. Carbon chemistry parameters in this region are not only under the direct impact of Taiwan Warm Current and remote control of the Kuroshio water intrusion into the shelf but also significantly controlled by seasonal variations of the Changjiang River discharge (e.g., Isobe and Matsuno, 2008; Chen et al., 2008; Chou et al., 2009). Taking into account the highly complex hydrographic, biological and chemical conditions, the accuracy of pH_T presented is promising.

145 **3.2 Comparison with new field data**

To further assess the ability of the ANN model to estimate pH_T on the ECS shelf, we applied the ANN model to data from three cruises not used in the ANN model development (Fig. 6): March, July, and October 2018. Scatterplots of retrieved pH_T vs observations (Fig. 7) show that the ANN model predicts pH_T with a RMSE of 0.04 and R^2 of 0.8 for these cruise data. This result is consistent with the result of the training data (Fig. 3a), which further reflects the stability and reliability of the ANN
150 model on the ECS shelf.

3.3 Variable importance in the ANN model

To assess the relative importance of different environmental input variables in the ANN model, we used the following method: for each environmental variable separately, add 5% and calculate the resulting percentage change in the predicted pH_T . Predicted pH_T responded positively to (T, DO) and negatively to S (Fig. 8). The variable with the greatest weight was DO,
155 followed by S and T, and the weights of nutrients were relatively small. This is consistent with (Cai et al., 2011) where positive correlations between pH_T and DO in the Gulf of Mexico and ECS were attributed to the processes of photosynthesis (generating DO and removing CO_2 , hence increasing pH) and aerobic respiration (consuming DO and generating CO_2 , hence lowering pH). The negative response to increasing S reflects the influence of the Changjiang River (lower salinity), which carries large



amounts of nutrients that fuel increased primary production (uptake of nutrients and CO₂, hence raising the pH) in surface
160 waters during warm seasons (Gong et al., 2011).

3.4 ANN model application

In order to retrieve pH_T on the ECS shelf, the monthly T, S, DO, N, P and Si from the Changjiang Biology Finite-Volume Coastal Ocean Model (FVCOM) (<http://47.101.49.44/wms/demo>) were applied to the ANN model as input variables.

165 Monthly pH_T for the period 2000-2016 was obtained at the spatial resolution of the Changjiang Biology FVCOM output: 1-10 km in the horizontal, 10 depth levels in the vertical, and 12 months.

Considering the discreteness and discontinuity of the sampling sites, we compared retrieved pH_T with the corresponding observations at some sites with repeated sampling for 3 to 4 years. These sites were A1-5 (123.0140 °E, 32.2145 °N), A1-6 (123.2750 °E, 32.2679 °N), A6-7 (122.9880 °E, 30.7050 °N), A6-9 (123.4990 °E, 30.5723 °N), A7-5 (123.4990 °E, 30.2523 °N), and A8-5 (123.4930 °E, 29.9940 °N). Overall, the retrieved pH_T from the Changjiang Biology FVCOM output agrees well with
170 the observed values at the surface, except for three samples in summer (Fig. 9). There are relatively large deviations (greater than the RMSE of 0.04) in August 2013 at station A1-5 and A6-9, and in July 2016 at station A8-5. These may be primarily attributed to the sudden increase in the Changjiang River discharge (Dai and Trenberth, 2002). To illustrate the application performance in the water column, a scatterplot of retrieved pH_T vs observations at six sites with repeated sampling for 3 to 4 years (Fig. 10) shows that the ANN model predicts pH_T with a RMSE of 0.05 and R² of 0.71.

175 We further compared retrieved pH_T using the Changjiang Biology FVCOM output with retrieved pH_T using measured T, S and DO, and in situ measured pH_T values (Figure 11). The agreement is good here in winter, but large deviations appear in summer. The reduced performance in summer can be attributed in large part a reduced performance of the FVCOM model in predicting summertime DO and S (see Figure S1); using the observed values of DO, S, etc. as predictor variables, the skill of the ANN pH_T predictions is much improved (RMSE = 0.09 vs. RMSE = 0.02).

180 4 Conclusions and perspectives

We have developed an artificial neural network model, demonstrated its reliability, and used it to retrieve monthly pH_T for the period 2000-2016 on the East China Sea shelf. This model predicts the water column pH_T using nine input components, and the three most important environmental input variables were dissolved oxygen, salinity and temperature.

185 The approach has several potential applications. First, it can provide estimates of seawater pH_T with known accuracies for the East China Sea shelf and the period 2013-2018. Within this region the model could be used as a cost-effective way to handle restrictions of marine observations conducted from ships, such as coarse resolution and under-sampling of carbonate system variables. Second, while the ANN model is not a replacement for direct measurements of the carbonate system, it may be a valuable tool for understanding the seasonal variation of pH_T in poorly observed regions. Third, this approach can be applied
190 to other regions to predict pH by suitably adapting the input variables and network structure. The MATLAB code used in this study to develop and apply the ANN model is freely available, and is accompanied by a README file providing detailed guidance on how to use and adapt the code.

Code and data availability

Matlab code of the ANN model for pH_T estimation and datasets are available:

195 <http://doi.org/10.5281/zenodo.3519219>



The monthly-average input variables (T, S, DO, N, P, Si) from the Changjiang Biology Finite-Volume Coastal Ocean Model and retrieved pH_T values from 2000 to 2016 on the East China Sea shelf and three cruises data during 2018 used to evaluate the ANN model are available:

<http://doi.org/10.5281/zenodo.3519236>

200 Requests to access the raw data should be directed to Richard Bellerby: Richard.Bellerby@niva.no

Six stations with repeated sampling for 3 to 4 years and corresponding retrieved pH values from the Changjiang Biology FVCOM output are available: <http://doi.org/10.5281/zenodo.3491747>

Video supplement

Monthly distribution of surface pH_T on the East China Sea shelf from 2000 to 2016 year:

205 <http://doi.org/10.5281/zenodo.2672943>

Profile distribution of pH_T at 31°N on the East China Sea shelf from 2000 to 2016 year:

<http://doi.org/10.5281/zenodo.2672929>

Author contribution

210 Li, X. S. and Bellerby, R. contributed to the development of methodology and the design of the model. Ge, J. Z. provided ten cruises dataset from 2013 to 2017 year and the input variables from the Changjiang Biology Finite-Volume Coastal Ocean Model Data. Liu, J. and Yang, A. Q. provided four cruises dataset from 2017 to 2018 year. Li, X. S. developed the manuscript with contributions from all co-authors.

Acknowledgements

215 This study was financially supported by the National Thousand Talents Program for Foreign Experts (grants No. WQ20133100150), Vulnerabilities and Opportunities of the Coastal Ocean (grants No. SKLEC-2016RCDW01), Marginal Seas (MARSEAS) (grants SKLEC-Taskteam project), and Innovative Talents International Cooperation Training Project (grants No. China Scholarship Council-201913045). Richard Bellerby and Philip Wallhead were also supported by funding from the FRAM High North Research Centre for Climate and the Environment under the Ocean Acidification Flagship and
220 the NIVA Land-Ocean Interactions Strategic Institute program. We deeply thank the people who worked on the cruises and in the laboratory.

References

- Alin, S. R., Feely, R. A., Dickson, A. G., Hernández-Ayón, J. M., Juranek, L. W., Ohman, M. D., and Goericke, R.: Robust empirical relationships for estimating the carbonate system in the southern California Current System and application to
225 CalCOFI hydrographic cruise data (2005–2011), *J. Geophys. Res.*, 117, C05033, doi:10.1029/2011JC007511, 2012.
- Bai, Y., Cai, W. J., He, X. Q., Zhai, W. D., Pan, D., Dai, M. H., and Yu, P. S.: A mechanistic semi-analytical method for remotely sensing sea surface pCO_2 in river-dominated coastal oceans: A case study from the East China Sea, *J. Geophys. Res. Oceans*, 120, 2331–2349, doi:10.1002/2014JC010632, 2015.
- Bates, N. R., Astor, Y. M., Church, M. J., Currie, K., Dore, J. E., González-Dávila, M., Lorenzoni, L., Muller-Karger, F.,
230 Olafsson, J., and Santana-Casiano, J. M.: A time-series view of changing ocean chemistry due to ocean uptake of anthropogenic CO_2 and ocean acidification, *Oceanography*, 27(1), 126–141, doi:10.5670/oceanog.2014.16, 2014.



- Bostock, H. C., Mikaloff Fletcher, S. E., and Williams, M. J. M.: Estimating carbonate parameters from hydrographic data for the intermediate and deep waters of the Southern Hemisphere oceans, *Biogeosciences*, 10, 6199–6213, <https://doi.org/10.5194/bg-10-6199-2013>, 2013.
- 235 Cai, W. J., Hu, X. P., Huang W. J., Murrell, M. C., Lehrter, J. C., Lohrenz, S. E., Chou, W. C., Zhai, W. D., Hollibaugh, J. T., Wang, Y. C., Zhao, P. S., Guo, X. H., Gundersen, K., Dai, M. H., and Gong, G. C.: Acidification of subsurface coastal waters enhanced by eutrophication, *Nature Geoscience*, 4, 766–770, doi:10.1038/NGEO1297, 2011.
- Cao, Z. M., Dai, M. H., Zheng, N., Wang, D., Li, Q., Zhai, W. D., Meng, F. F., and Gan, J. P.: Dynamics of the carbonate system in a large continental shelf system under the influence of both a river plume and coastal upwelling, *J. Geophys. Res.*, 240 116, G02010, doi:10.1029/2010JG001596, 2011.
- Carter, B. R., Feely, R. A., Williams, N. L., Dickson, A. G., Fong, M. B., and Takeshita, Y.: Updated methods for global locally interpolated estimation of alkalinity, pH, and nitrate, *Limnol. Oceanogr. Methods*, 16, 119–131, doi:10.1002/lom3.10232, 2018.
- Carter, B. R., Williams, N. L., Gray, A. R., and Feely, R. A.: Locally interpolated alkalinity regression for global alkalinity estimation, *Limnol. Oceanogr. Methods*, 14, 268–277, doi:10.1002/lom3.10087, 2016.
- 245 Chen, C. S., Xue, P. F., Ding, P. X., Beardsley, R. C., Xu, Q. C., Mao, X. M., Gao, G. P., Qi, J. H., Li, C. Y., Lin, H. C., Cowles, G., and Shi, M. C.: Physical mechanisms for the offshore detachment of the Changjiang Diluted Water in the East China Sea, *J. Geophys. Res.*, 113, C02002, doi:10.1029/2006JC003994, 2008.
- Chen, S. L. and Hu, C. M.: Estimating sea surface salinity in the northern Gulf of Mexico from satellite ocean color measurements, *Remote Sens Environ*, 201, 115–132, <https://doi.org/10.1016/j.rse.2017.09.004>, 2017.
- 250 Chou, W. C., Gong, G. C., Sheu, D. D., Hung, C. C., and Tseng, T. F.: Surface distributions of carbon chemistry parameters in the East China Sea in summer 2007, *J. Geophys. Res.*, 114, C07026, doi:10.1029/2008JC005128, 2009.
- Ciais, P., Sabine, C., Bala, G., Bopp, L., Brovkin, V., Canadell, J., Chhabra, A., DeFries, R., Galloway, J., Heimann, M., Jones, C., Le Quéré C., Myneni, R. B., Piao, S., and Thornton, P.: Carbon and Other Biogeochemical Cycles. In: *Climate Change 2013: The Physical Science Basis. Contribution of Working Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate Change* [Stocker, T. F., D. Qin, G. K. Plattner, M. Tignor, S. K. Allen, J. Boschung, A. Nauels, Y. Xia, V. Bex and P. M. Midgley (eds.)]. Cambridge University Press, Cambridge, United Kingdom and New York, NY, USA, 2013.
- 255 Dai, A. and Trenberth, K. E.: Estimates of freshwater discharge from continents: Latitudinal and seasonal variations, *J. Hydrometeorol.*, 3, 660–687, <https://doi.org/10.1175/1525-7541>, 2002.
- 260 Dore, J., Lukas, R., Sadler, D., Church, M., and Karl, D.: Physical and biogeochemical modulation of ocean acidification in the central North Pacific, *Proc. Natl. Acad. Sci. U. S. A.*, 106, 12235–12240, 2009.
- Friedrich, T. and Oschlies, A.: Neural network-based estimates of North Atlantic surface pCO₂ from satellite data: A methodological study, *J. Geophys. Res.*, 114, C03020, <https://doi.org/10.1029/2007JC004646>, 2009.
- Gong, G. C., Chen, Y. L. L., and Liu, K. K.: Chemical hydrography and chlorophyll a distribution in the East China Sea in summer: implications in nutrient dynamics, *Cont. Shelf Res*, 16, 1561–1590, [https://doi.org/10.1016/0278-4343\(96\)00005-2](https://doi.org/10.1016/0278-4343(96)00005-2), 1996.
- 265 Gong, G. C., Liu, K. K., Chiang, K. P., Hsiung, T. M., Chang, J., Chen, C. C., Hung, C. C., Chou, W. C., Chung, C. C., Chen, H. Y., Shiah, F. K., Tsai, A. Y., Hsieh, C. H., Shiao, J. C., Tseng, C. M., Hsu, S. C., Lee, H. J., Lee, M. A., Lin, I. I., and Tsai, F.: Yangtze River floods enhance coastal ocean phytoplankton biomass and potential fish production, *Geophys. Res. Lett.*, 38, L13603, doi:10.1029/2011GL047519, 2011.
- 270 González-Dávila, M., Santana-Casiano, J. M., Rueda, M. J., and Llinás, O.: The water column distribution of carbonate system variables at the ESTOC site from 1995 to 2004, *Biogeosciences*, 7, 3067–3081, 2010.
- Guo, X. H., Zhai, W. D., Dai, M. H., Zhang, C., Bai, Y., Xu, Y., Li, Q., and Wang, G. Z.: Air–sea CO₂ fluxes in the East China Sea based on multiple-year underway observations, *Biogeosciences*, 12, 5495–5514, doi:10.5194/bg-12-5495-2015, 2015.



- 275 Hornik, K., Stinchcombe, M., and White, H.: Multilayer feedforward networks are universal approximators, *Neural Netw.*, 2, 359–366, [https://doi.org/10.1016/0893-6080\(89\)90020-8](https://doi.org/10.1016/0893-6080(89)90020-8), 1989.
- Isobe, A. and Matsuno, T.: Long-distance nutrient-transport process in the Changjiang River plume on the East China Sea shelf in summer, *J. Geophys. Res.-Oceans*, 113, C04006, doi:10.1029/2007JC004248, 2008.
- Juranek, L. W., Feely, R. A., Gilbert, D., Freeland, H., and Miller, L. A.: Real-time estimation of pH and aragonite saturation state from Argo profiling floats: Prospects for an autonomous carbon observing strategy, *Geophys. Res. Lett.*, 38, L17603, doi:10.1029/2011GL048580, 2011.
- 280 Laruelle, G. G., Landschützer, P., Gruber, N., Tison, J. L., Delille, B., and Regnier, P.: Global high-resolution monthly pCO₂ climatology for the coastal ocean derived from neural network interpolation, *Biogeosciences*, 14, 4545–4561, doi:10.5194/bg-14-4545-2017, 2017.
- 285 Lauvset, S. K., Gruber, N., Landschützer, P., Olsen, A., and Tjiputra, J.: Trends and drivers in global surface ocean pH over the past 3 decades, *Biogeosciences*, 12(5), 1285–1298, doi:10.5194/bg-12-1285-2015, 2015.
- Le Quéré, C., Andrew, R. M., Friedlingstein, P., Sitch, S., Hauck, J., Pongratz, J., and Pickers, P. A., et al.: Global Carbon Budget 2018, *Earth Syst. Sci. Data*, 10, 2141–2194, <https://doi.org/10.5194/essd-10-2141-2018>, 2018.
- Lueker, T. J., Dickson, A. G., and Keeling, C. D.: Ocean pCO₂ calculated from dissolved inorganic carbon, alkalinity, and equations for K₁ and K₂: Validation based on laboratory measurements of CO₂ in gas and seawater at equilibrium, *Mar. Chem.*, 290, 105–119, doi:10.1016/S0304-4203(00)00022-0, 2000.
- Moore-Maley, B. L., Allen, S. E., and Ianson, D.: Locally driven interannual variability of near-surface pH and Ω_A in the Strait of Georgia, *J. Geophys. Res. Oceans*, 121, 1600–1625, doi:10.1002/2015JC011118, 2016.
- Olden, J. D. and Jackson, D. A.: Illuminating the “black box”: a randomization approach for understanding variable contributions in artificial neural networks, *Ecol. Model.*, 154, 135–150, [https://doi.org/10.1016/S0304-3800\(02\)00064-9](https://doi.org/10.1016/S0304-3800(02)00064-9), 2002.
- 295 Olden, J. D., Joy, M. K., and Death, R. G.: An accurate comparison of methods for quantifying variable importance in artificial neural networks using simulated data, *Ecol. Model.*, 178, 389–397, <https://doi.org/10.1016/j.ecolmodel.2004.03.013>, 2004.
- Palacz, A. P., John, M. A. S., Brewin, R. J. W., Hirata, T., and Gregg, W. W.: Distribution of phytoplankton functional types in high-nitrate, low-chlorophyll waters in a new diagnostic ecological indicator model, *Biogeosciences*, 10, 8103–8157, 300 <https://doi.org/10.5194/bg-10-8103-2013>, 2013.
- Pelletier, G. J., Lewis, E., and Wallace, D. W. R.: CO₂SYX.XLS: A calculator for the CO₂ system in seawater for Microsoft Excel/VBA, Ver. 16, Washington State Department of Ecology, Olympia, Washington, 2011.
- Qu, B. X., Song, J. M., Yuan, H. M., Li, X. G., Li, N., Duan, L. Q., Chen, X., and Lu, X.: Summer carbonate chemistry dynamics in the Southern Yellow Sea and the East China Sea: Regional variations and controls, *Cont. Shelf Res.*, 111, 250–261, <https://doi.org/10.1016/j.csr.2015.08.017>, 2015.
- 305 Raitos, D. E., Lavender, S. J., Maravelias, C. D., Haralabous, J., Richardson, A. J., and Reid, P. C.: Identifying four phytoplankton functional types from space: an ecological approach, *Limnol. Oceanogr.*, 53, 605–613, <https://doi.org/10.4319/lo.2008.53.2.0605>, 2008.
- Reggiani, E. R., King, A. L., Norli, M., Jaccard, P., Sorensen, K., and Bellerby, R. G. J.: FerryBox-assisted monitoring of mixed layer pH in the Norwegian Coastal Current, *Journal of Marine Systems*, 162, 29–36, doi:10.1016/j.jmarsys.2016.03.017, 2016.
- 310 Sabine, C. L., Feely, R. A., Gruber, N., Key, R. M., Lee, K., Bullister, J. L., Wanninkhof, R., Wong, C. S., Wallace, D. W. R., Tilbrook, B., Millero, F. J., Peng, T. H., Kozyr, A., Ono, T., Rios, A. F.: The oceanic sink for anthropogenic CO₂, *Science*, 305, 367–371, 2004.
- 315 Sasse, T. P., McNeil, B. I., and Abramowitz, G.: A novel method for diagnosing seasonal to inter-annual surface ocean carbon dynamics from bottle data using neural networks, *Biogeosciences*, 10, 4319–4340, doi:10.5194/bg-10-4319-2013, 2013.



- Sauzède, R., Bittig, H. C., Claustre, H., de Fommervault, O. P., Gattuso, J. P., Legendre, L., and Johnson, K. S.: Estimates of Water-Column Nutrient Concentrations and Carbonate System Parameters in the Global Ocean: A Novel Approach Based on Neural Networks, *Front. Mar. Sci.*, 4, 128, doi:10.3389/fmars.2017.00128, 2017.
- 320 Sauzède, R., Claustre, H., Jamet, C., Uitz, J., Ras, J., Mignot, A., and D'Ortenzio, F.: Retrieving the vertical distribution of chlorophyll-a concentration and phytoplankton community composition from in situ fluorescence profiles: a method based on a neural network with potential for global-scale applications, *J. Geophys. Res. Ocean*, 120, 451–470, doi:10.1002/2014JC010355, 2015.
- Sauzède, R., Claustre, H., Uitz, J., Jamet, C., Dall'Olmo, G., D'Ortenzio, F., Gentili, B., Poteau, A., and Schmechtig, C.: A neural network-based method for merging ocean color and Argo data to extend surface bio-optical properties to depth: retrieval of the particulate backscattering coefficient, *J. Geophys. Res. Ocean*, 121, 2552–2571, doi:10.1002/2015JC011408, 2016.
- 325 Shim, J. H., Kim, D., Kang, Y. C., Lee, J. H., Jang, S. T., and Kim, C. H.: Seasonal variations in pCO₂ and its controlling factors in surface seawater of the northern East China Sea, *Cont. Shelf Res.*, 27, 2623–2636, <https://doi.org/10.1016/j.csr.2007.07.005>, 2007.
- 330 Tamura, S. and Tateishi, M.: Capabilities of a Four-Layered Feedforward Neural Network: Four Layers versus Three, *IEEE Transactions on Neural Networks*, 8, 251–255, doi:10.1109/72.557662, 1997.
- Uusitalo, L.: Advantages and challenges of Bayesian networks in environmental modelling, *Ecol. Model*, 203, 312–318, <https://doi.org/10.1016/j.ecolmodel.2006.11.033>, 2007.
- Velo, A., Pérez, F. F., Tanhua, T., Gilcoto, M., Ríos, A. F., and Key, R. M.: Total alkalinity estimation using MLR and neural network techniques, *J Marine Syst*, 111–112, 11–18, <https://doi.org/10.1016/j.jmarsys.2012.09.002>, 2013.
- 335 Williams, N. L., Juranek, L. W., Johnson, K. S., Feely, R. A., Riser, S. C., Talley, L. D., Russell, J. L., Sarmiento, J. L., and Wanninkhof, R.: Empirical algorithms to estimate water column pH in the Southern Ocean, *Geophys. Res. Lett.*, 43, 3415–3422, doi:10.1002/2016GL068539, 2016.
- Wootton, J. T., Pfister, C. A., and Forester, J. D.: Dynamic patterns and ecological impacts of declining ocean pH in a high resolution multi-year dataset, *Proc. Natl. Acad. Sci. U. S. A.*, 105, 18848–18853, <https://doi.org/10.1073/pnas.0810079105>, 2008.
- 340 Wootton, J. T. and Pfister, C. A.: Carbon System Measurements and Potential Climatic Drivers at a Site of Rapidly Declining Ocean pH, *PLoS ONE*, 7(12): e53396, <https://doi.org/10.1371/journal.pone.0053396>, 2012.
- Zhai, W. D. and Dai, M. H.: On the seasonal variation of air-sea CO₂ fluxes in the outer Changjiang (Yangtze River) Estuary, East China Sea, *Mar. Chem.*, 117, 2–10, <https://doi.org/10.1016/j.marchem.2009.02.008>, 2009.
- 345 Zhai, W. D., Zhao, H. D., Zheng, N., and Xu, Y.: Coastal acidification in summer bottom oxygen-depleted waters in northwestern-northern Bohai Sea from June to August in 2011, *Chinese Science Bulletin*, 57, 1062–1068, doi:10.1007/s11434-011-4949-2, 2012.
- Zhai, W. D., Zheng, N., Huo, C., Xu, Y., Zhao, H. D., Li, Y. W., Zang, K. P., Wang, J. Y., and Xu, X. M.: Subsurface pH and carbonate saturation state of aragonite on the Chinese side of the North Yellow Sea: seasonal variations and controls, *Biogeosciences*, 11, 1103–1123, <https://doi.org/10.5194/bg-11-1103-2014>, 2014.
- 350 Zhang, G., Zhang, J., and Liu, S. M.: Characterization of nutrients in the atmospheric wet and dry deposition observed at the two monitoring sites over Yellow Sea and East China Sea, *J Atmos Chem.*, 57(1), 41–57, doi:10.1007/s10874-007-9060-3, 2007.
- 355

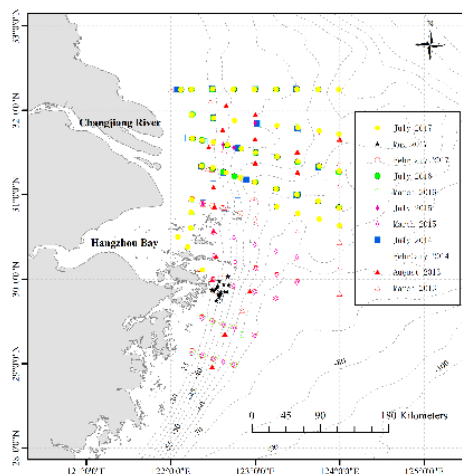


Figure 1: Sampling stations from 11 cruises from 2013 to 2017.

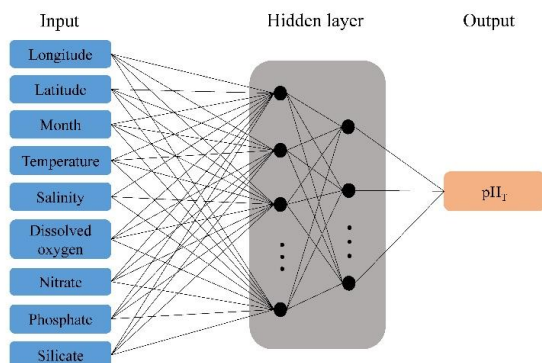
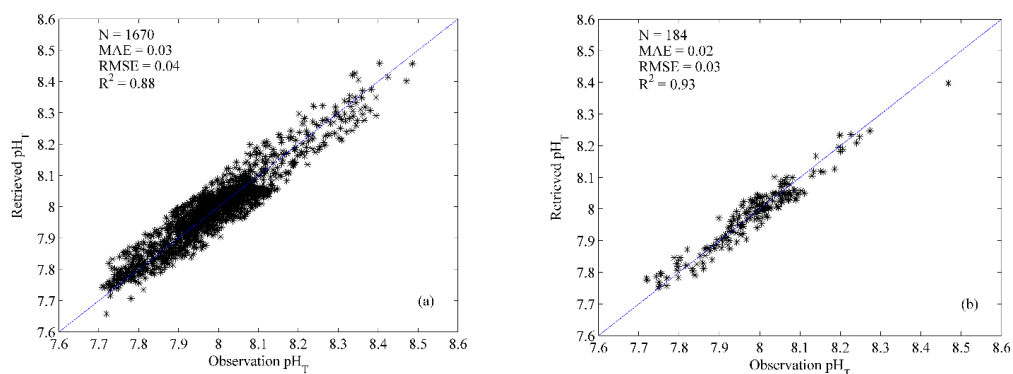
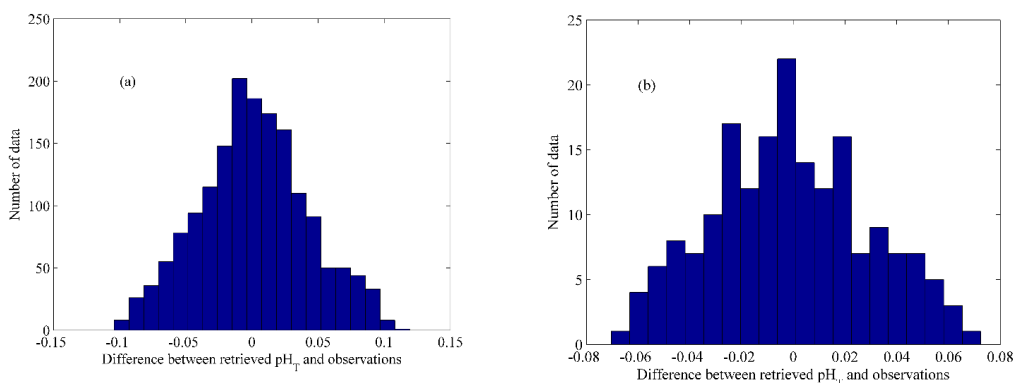


Figure 2: Schematic representation of the neural network algorithm to retrieve pH_T . Input variables are observed temperature, salinity, dissolved oxygen, nitrate, phosphate, and silicate together with the geolocation (longitude and latitude) and time (month) of sampling.



360 Figure 3: Comparison of pH_T retrieved by the ANN model with corresponding observations. (a)-training data (90% of all data); (b)-independent validation data (10% of all data). The 1:1 line is shown in each plot as visual reference. Three statistics are the mean absolute error (MAE), the coefficient of determination (R^2), and the root mean squared error (RMSE). N represents the number of data points.



365 **Figure 4: Normal distribution of the differences between pH_T estimated by the ANN model minus the observations. (a)-training data (90% of all data); (b)-independent validation data (10% of all data).**

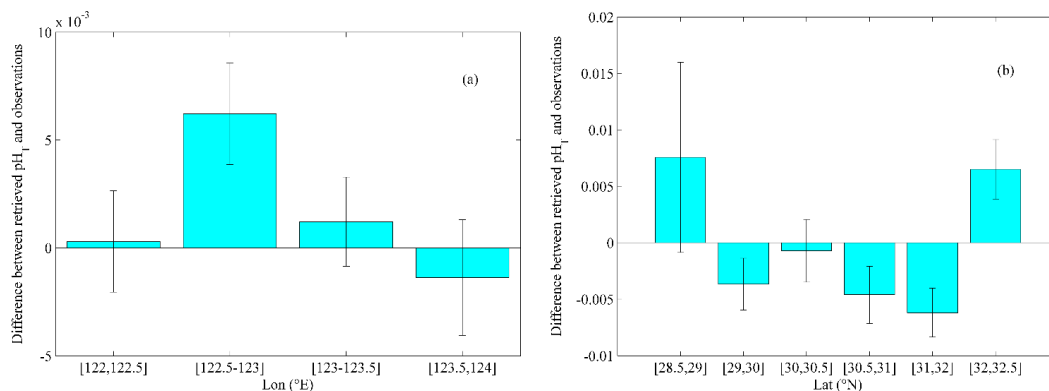
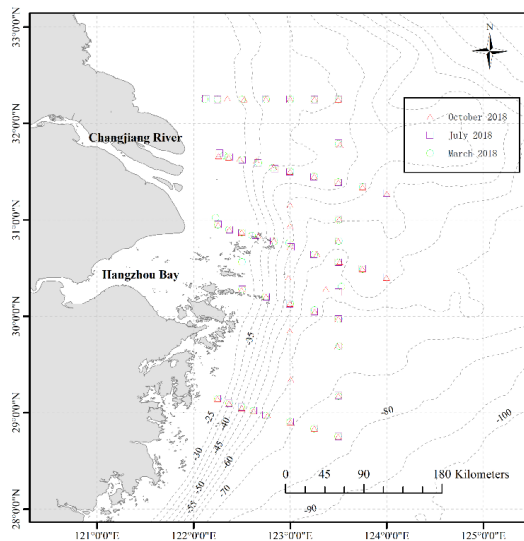


Figure 5: Box plots of the differences between retrieved pH_T minus the observations. (a)-the differences vs longitude (Mean \pm SE); (b)-the differences vs latitude (Mean \pm SE). The height of each box represents the mean value of the differences, the whisker represents the standard error (SE) value of the differences.



370 **Figure 6: Sampling stations for three cruises used to extend the utility of the ANN model. The green circles represent March 2018, the purple squares represent July 2018, the red triangles represent October 2018.**

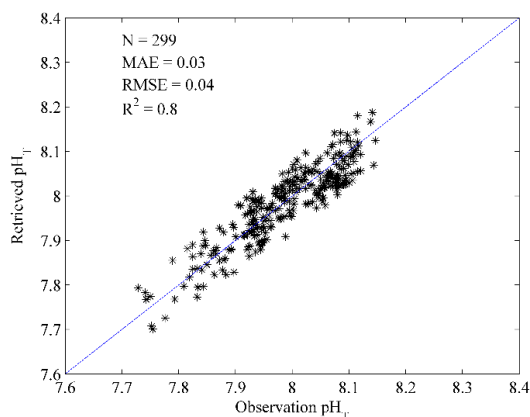


Figure 7: Comparison of pH_T retrieved by the ANN model with corresponding observations for the three cruises in 2018. The 1:1 line is shown in the plot as visual reference. Three statistics approaches used are the mean absolute error (MAE), the coefficient of determination (R^2), and the root mean squared error (RMSE). N represents the number of data points.

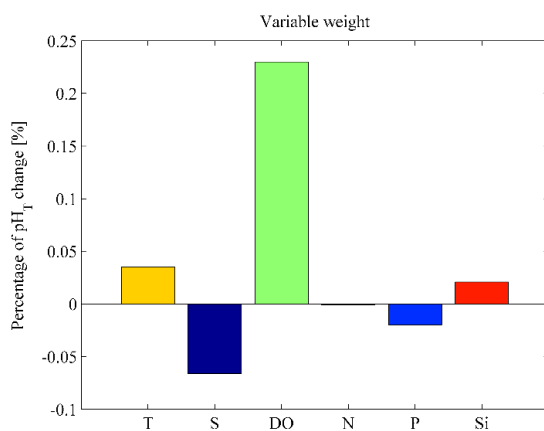
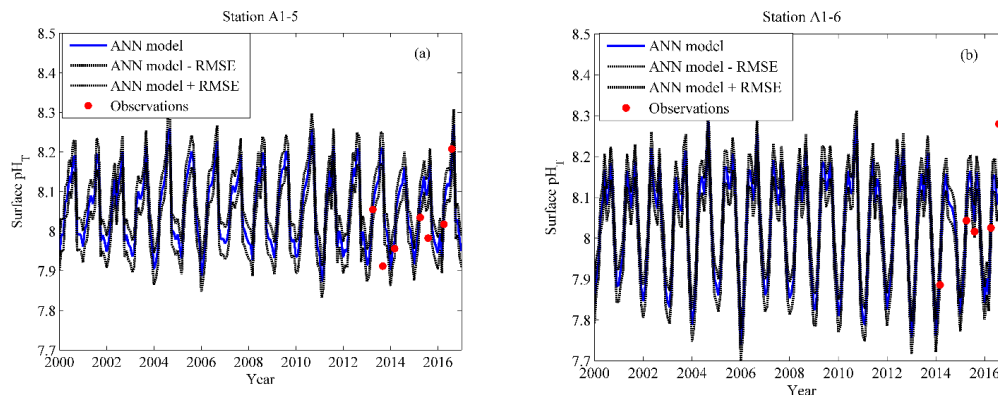
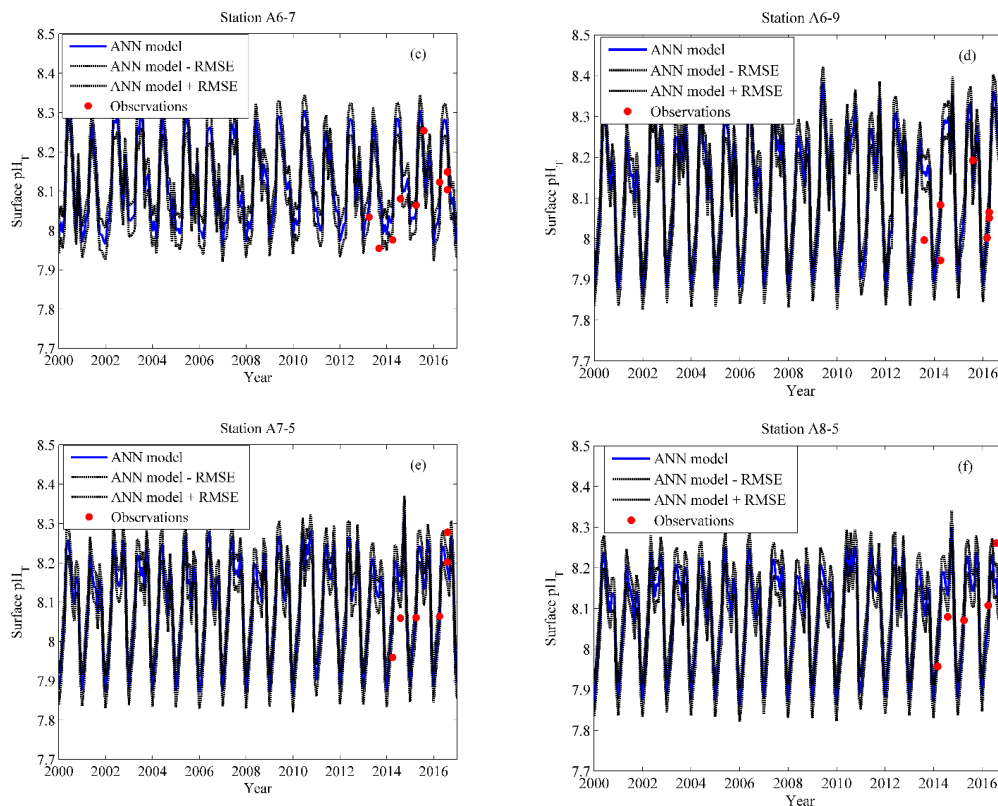
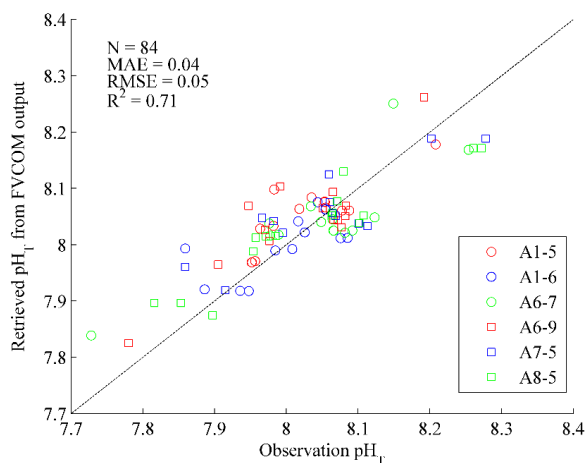


Figure 8: Variable importance estimates in the ANN model.

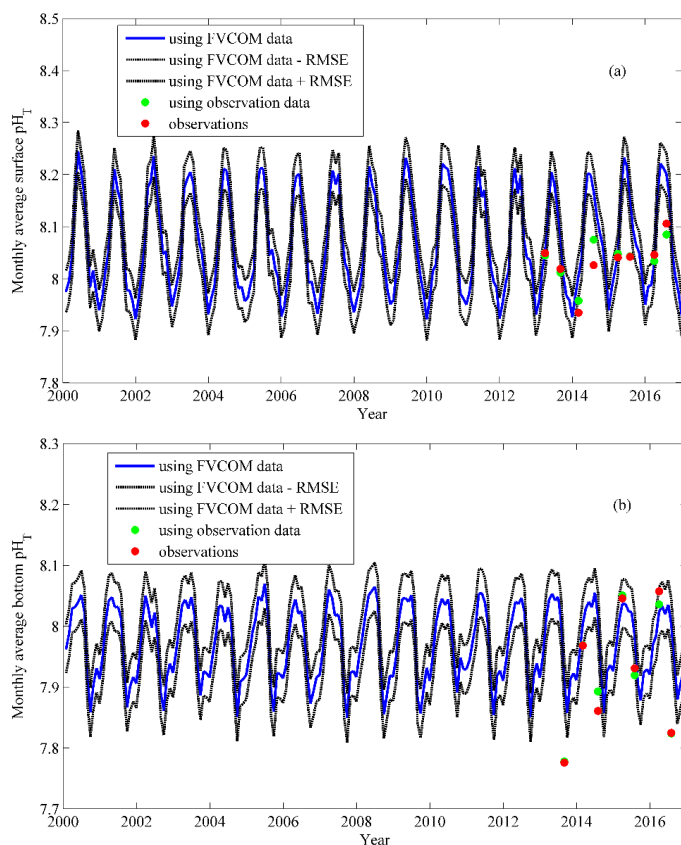




375 **Figure 9: Comparison of retrieved pH_T with corresponding observations in the sea surface at six sites repeated sampling for 3 to 4 years. Red dots represent observations pH_T , blue solid line represents time series pH_T retrieved by the ANN model, black dotted line represent retrieved $pH_T \pm RMSE$. (a)-station A1-5; (b)-station A1-6; (c)-station A6-7; (d)-station A6-9; (e)-station A7-5; (f)-station A8-5.**



380 **Figure 10: Comparison of retrieved pH_T with corresponding observations in the water column at six sites repeated sampling for 3 to 4 years. The 1:1 line is shown in the plot as a visual reference. Skill statistics include the mean absolute error (MAE), the coefficient of determination (R^2), and the root mean squared error (RMSE). N represents the number of data points.**



385 **Figure 11: Comparison of monthly average pH_T on the East China Sea shelf. Blue solid line represents retrieved pH_T using Changjiang Biology FVCOM output; black dotted line represents retrieved pH_T using Changjiang Biology FVCOM output \pm RMSE; red points show monthly-average pH_T observations from 2013 to 2016; green points show retrieved pH_T using in situ observation data from 2013 to 2016. (a)-surface; (b)-bottom.**

Table 1: Field survey information and measurements of water temperature, salinity, dissolved oxygen, nitrate, phosphate, silicate and pH_T (Mean \pm SE).

Sampling period	Temperature (°C)	Salinity	Dissolved oxygen (mmol m ⁻³)	Nitrate (mmol m ⁻³)	Phosphate (mmol m ⁻³)	Silicate (mmol m ⁻³)	pH_T
March 4 th -20 th , 2013	11.54 \pm 1.34	32.04 \pm 2.26	275.28 \pm 19.30	12.25 \pm 8.25	0.58 \pm 0.17	17.54 \pm 7.65	8.19 \pm 0.04
August 17 th -28 th , 2013	23.45 \pm 3.17	32.32 \pm 2.91	142.22 \pm 63.45	12.16 \pm 8.05	0.55 \pm 0.32	16.47 \pm 12.18	8.04 \pm 0.18
February 21 th -28 th , 2014	9.56 \pm 2.38	32.14 \pm 1.78	293.07 \pm 19.52	11.92 \pm 9.17	0.59 \pm 0.18	12.52 \pm 6.50	8.10 \pm 0.04
July 10 th -17 th , 2014	21.66 \pm 2.13	29.50 \pm 5.10	186.44 \pm 43.29	21.57 \pm 22.10	0.57 \pm 0.46	21.45 \pm 17.76	8.07 \pm 0.11
March 11 th -21 th , 2015	11.42 \pm 1.44	31.57 \pm 2.60	279.72 \pm 15.29	22.04 \pm 18.88	0.81 \pm 0.35	16.48 \pm 11.64	8.19 \pm 0.03
July 9 th -20 th , 2015	22.14 \pm 1.55	29.73 \pm 4.71	207.32 \pm 56.12	19.73 \pm 18.62	0.60 \pm 0.42	20.87 \pm 17.48	8.13 \pm 0.09
March 7 th -19 th , 2016	10.77 \pm 2.02	30.85 \pm 2.92	284.00 \pm 31.40	20.26 \pm 12.80	0.82 \pm 0.25	19.17 \pm 11.62	8.20 \pm 0.05
July 4 th -28 th , 2016	23.19 \pm 3.19	28.17 \pm 6.67	122.90 \pm 49.97	25.77 \pm 23.60	0.63 \pm 0.46	28.56 \pm 25.03	8.06 \pm 0.16
February 15 th -28 th , 2017	11.03 \pm 2.57	32.00 \pm 2.43	296.21 \pm 21.27	12.30 \pm 9.13	0.56 \pm 0.18	13.09 \pm 7.45	8.13 \pm 0.05
May 12 th -24 th , 2017	17.71 \pm 1.54	29.62 \pm 2.79	171.58 \pm 49.52	12.60 \pm 4.83	0.29 \pm 0.24	10.95 \pm 4.29	8.08 \pm 0.13
July 20 th -30 th , 2017	24.85 \pm 3.41	27.70 \pm 6.31	192.11 \pm 76.55	20.57 \pm 23.23	0.42 \pm 0.34	19.28 \pm 18.92	8.09 \pm 0.18

390

Table 2: Different model structures and their performance in the training step. The variables (Lon (longitude), Lat (latitude), Month (month), T (temperature), S (salinity), DO (dissolved oxygen), N (nitrate), P (phosphate), Si (silicate)) marked with 1 represent the input variables. Skill statistics include the coefficient of determination (R^2), the root mean squared error (RMSE), and the mean absolute error (MAE).

395



Model	Lon	Lat	Month	T	S	DO	N	P	Si	Training data			Independent validation data			
										R ²	RMSE	MAE	R2	RMSE	MAE	
1						1					0.40	0.092	0.068	0.47	0.076	0.058
2				1		1					0.62	0.073	0.053	0.62	0.067	0.051
3				1	1	1					0.69	0.065	0.048	0.72	0.060	0.044
4				1	1	1	1				0.76	0.057	0.044	0.77	0.052	0.041
5				1	1	1		1			0.81	0.051	0.040	0.79	0.051	0.040
6				1	1	1			1		0.77	0.056	0.044	0.79	0.054	0.043
7				1	1	1	1	1			0.80	0.053	0.042	0.79	0.051	0.041
8				1	1	1		1	1		0.81	0.051	0.040	0.81	0.049	0.039
9				1	1	1	1		1		0.76	0.058	0.044	0.77	0.054	0.044
10				1	1	1	1	1	1		0.83	0.048	0.037	0.86	0.046	0.037
11			1	1	1	1	1	1			0.85	0.046	0.035	0.87	0.043	0.032
12			1	1	1	1		1	1		0.85	0.046	0.034	0.85	0.045	0.035
13			1	1	1	1	1		1		0.82	0.049	0.036	0.84	0.050	0.036
14			1	1	1	1	1	1	1		0.84	0.046	0.035	0.87	0.045	0.033
15	1	1	1	1	1	1	1				0.86	0.044	0.033	0.79	0.046	0.034
16	1	1	1	1	1	1		1			0.87	0.043	0.032	0.87	0.044	0.034
17	1	1	1	1	1	1			1		0.87	0.043	0.033	0.82	0.045	0.035
18	1	1	1	1	1	1	1	1			0.88	0.040	0.031	0.88	0.039	0.031
19	1	1	1	1	1	1		1	1		0.87	0.042	0.032	0.87	0.042	0.033
20	1	1	1	1	1	1	1		1		0.84	0.046	0.035	0.85	0.047	0.036
21	1	1	1	1	1	1	1	1	1		0.88	0.040	0.031	0.89	0.036	0.028