



1 **Spatial reconstruction of long-term (2003-2020) sea surface $p\text{CO}_2$ in the**
2 **South China Sea using a machine learning based regression method**
3 **aided by empirical orthogonal function analysis**

4 Zhixuan Wang¹, Guizhi Wang^{1,2}, Xianghui Guo¹, Yan Bai³, Yi Xu¹ and Minhan Dai^{1,*}

5 ¹State Key Laboratory of Marine Environmental Science and College of Ocean and Earth Sciences, Xiamen University, Xiamen,
6 361102, China

7 ²Fujian Provincial Key Laboratory for Coastal Ecology and Environmental Studies, Xiamen University, Xiamen, 361102, China

8 ³State Key Laboratory of Satellite Ocean Environment Dynamics, Second Institute of Oceanography, State Oceanic
9 Administration, Hangzhou, 310012, China

10 *Correspondence to:* Minhan Dai (mdai@xmu.edu.cn)

11

12 **Abstract.** The South China Sea (SCS), the largest marginal sea of the North Pacific Ocean, is one of the world's most studied
13 model ocean margins in terms of its carbon cycle, where intensive field observations including sea-surface carbon dioxide partial
14 pressure ($p\text{CO}_2$) have been conducted over the last two decades. However, the datasets of cruise-based sea surface $p\text{CO}_2$ are still
15 temporally and spatially incomplete. Using a machine learning-based method facilitated by empirical orthogonal function (EOF)
16 analysis capable of constraining the spatiality, this study provides a reconstructed dataset of the monthly sea surface $p\text{CO}_2$ in the
17 SCS with a reasonably high spatial resolution ($0.05^\circ \times 0.05^\circ$) and temporal coverage between 2003 and 2020. We validate our
18 reconstruction with three independent testing datasets where, TEST.1 includes 10% of our observed data, TEST.2 includes four
19 independent underway datasets corresponding to four seasons, and TEST.3 includes a continuous observed dataset from 2003 -
20 2019 at the South East Asia Time-Series (SEATs) station located in the northern basin of the SCS. Our TEST.1 validation
21 demonstrated that the reconstructed $p\text{CO}_2$ field successfully simulated the spatial and temporal patterns of sea surface $p\text{CO}_2$. The
22 root-mean-square error (RMSE) between our reconstructions and observed data in TEST.1 averaged to $\sim 10 \mu\text{atm}$, which is much
23 smaller (by $\sim 50\%$) than that between the remote sensing (RS) and observed data. TEST.2 verified the accuracy of our
24 reconstruction model in data months lacking observations, showing a near-zero bias (RMSE: $\sim 8 \mu\text{atm}$). TEST.3 tested the
25 accuracy of the reconstructed long-term trend, showing that at the SEATs Station, the difference between the reconstructed $p\text{CO}_2$
26 and observations ranged from -10 to $4 \mu\text{atm}$ (-2.5 to 1%). In addition to the typical machine learning performance metrics, we
27 present a new method to assess the uncertainty that includes the bias from the reconstruction and its sensitivity to the features, and
28 successfully quantifies the spatial distribution patterns of uncertainty. These validations and uncertainty analysis strongly suggest
29 that our reconstruction is effectively captures the main features of both the spatial and temporal patterns of sea surface $p\text{CO}_2$ in the



30 SCS. Using the reconstructed dataset, we show the long-term trends of sea surface $p\text{CO}_2$ in 5 sub-regions of the SCS with
31 differing physico-biogeochemical characteristics. We show that mesoscale processes such as the Pearl River plume and China
32 Coastal Currents significantly impact sea surface $p\text{CO}_2$ in the SCS during different seasons. While the SCS is overall a weak
33 source of atmospheric CO_2 , the northern SCS acts as a sink, showing a trend of increasing strength over the past two decades.

34 Key words: Sea surface $p\text{CO}_2$; reconstruction; machine learning; South China Sea

35

36 **1 Introduction**

37 The ocean possesses much of the global capacity for atmospheric carbon dioxide (CO_2) sequestration and annually mitigates
38 22–26% of the anthropogenic CO_2 emissions associated with fossil fuel burning and land use change during the period 1960–2019
39 (Friedlingstein et al., 2020). However, it remains largely unknown whether and by how much the ocean will continue to act as a
40 sink for anthropogenic CO_2 , i.e., the extent of its climate change mitigation capacity to understand climate-carbon coupled
41 systems and develop zero-emission strategies and actions. Ocean margins, an essential part of the land-ocean continuum,
42 contributed $\sim 10\text{--}20\%$ of the global ocean CO_2 sequestration with only 7% of the surface area and represent a particularly
43 challenging regime (e.g., Chen and Borges, 2009; Dai et al. 2022; Laruelle et al., 2014). This is primarily attributed to the ocean
44 margins' extremely complex and dynamic processes, often characterized by large spatial and temporal variability of air-sea CO_2
45 fluxes that lead to even larger uncertainty in their prediction than those occurring in the open ocean (Dai et al., 2013, 2022; Cao et
46 al., 2020; Laruelle et al., 2014; Chen and Borges, 2009 and the references therein). Limited spatiotemporal coverage of
47 observational data is an important source of these uncertainties.

48 In recent years, many studies use numerical models or data-based approaches to improve estimates of sea surface CO_2 distribution
49 and the accuracy of the global carbon budget for periods and regions with poor coverage of observational data (Rödenbeck et al.,
50 2015; Wanninkhof et al., 2013). Numerical ocean models of performance can successfully quantify the generally increasing trend
51 in oceanic CO_2 and some critical processes of carbon cycling (e.g., net ecosystem production), but still suffer from regional and
52 seasonal differences in their estimates of the ocean carbonate system (Luo et al., 2015; Mongwe et al., 2016; Tahata et al., 2015;
53 Wanninkhof et al., 2013). Thus, data-based approaches have become a popular alternative to biogeochemical models (Jones et al.,
54 2014; Lefèvre et al., 2005; Landschützer et al., 2014, 2017; Telszewski et al., 2009). The former typically use statistical
55 interpolations and regression methods. Statistical interpolations improve the spatial coverage of observational data, but do not
56 work for the period without observational data. Regression methods allow mapping of the relationship between the observed
57 carbon dioxide partial pressure ($p\text{CO}_2$) data and other parameters that may drive changes in surface ocean $p\text{CO}_2$, and then
58 extrapolation of this relationship to improve estimates of the spatiotemporal distribution of $p\text{CO}_2$. The development of machine
59 learning methods and remote sensing-derived products (as proxy variables in regression methods) have aided the development of
60 data-based methods (Rödenbeck et al., 2015; Bakker et al., 2016) which, with spatial-temporal standardization, can improve the



61 model results of the oceanic carbonate system by numerical assimilation methods. However, because of the complex and dynamic
62 nature of biogeochemical and physical processes in coastal areas, characterization of sea surface $p\text{CO}_2$ and subsequently the
63 air-sea CO_2 fluxes both in time and space in marginal seas remains challenging.

64 The South China Sea (SCS) is the largest marginal sea of the North Pacific Ocean with a surface area of $3.5 \times 10^6 \text{ km}^2$. Although
65 extensive field observations have been conducted of sea surface $p\text{CO}_2$ in the SCS in the past two decades, their spatial and
66 temporal coverage is still limited in different physical-biogeochemical domains of the SCS and at sub-seasonal time scales (e.g.,
67 Guo et al., 2015; Li et al., 2020; Zhai et al., 2005; Zhai et al., 2013). Therefore, there is a clear need to achieve surface water $p\text{CO}_2$
68 coverage in the SCS with a highest spatiotemporal resolution as possible with the aim to better estimate sea surface $p\text{CO}_2$ and thus
69 air-sea CO_2 fluxes in the SCS and help develop improved initial conditions of numerical models. Moreover, a reasonably high
70 spatiotemporal resolution of $p\text{CO}_2$ data can help identify the controlling factors of $p\text{CO}_2$ changes in the SCS, and reliably resolve
71 the long-term changes.

72 Zhu et al. (2009) presented an empirical approach that estimated sea surface $p\text{CO}_2$ in the northern SCS in summer using
73 satellite-derived data (sea surface temperature, SST; chlorophyll a , Chl a), and their validation results show that the
74 reconstructed $p\text{CO}_2$ data was generally consistent with the underway observed data. However, it should be noted that the large
75 uncertainty of estimates from their study was caused by the limited underway observed data from only two summer cruises (July
76 2000 data were used for algorithm tuning and those of July 2000 for validation). Jo et al. (2012) developed a neural
77 networking-based algorithm by using SST and Chl a to estimate sea surface $p\text{CO}_2$ in the northern SCS. Sea surface $p\text{CO}_2$ data in
78 this study were collected from May 2001, and, February, July 2004. The difference between the reconstructed $p\text{CO}_2$ data of Jo et
79 al. (2012) and the observed data reflects a relatively large bias (the resultant RMSE (root-mean-square error) falls in the range
80 32.6 to 44.5 μatm , reported in Wang et al., 2021). Bai et al. (2015) used a ‘mechanic semi-analytical algorithm’ to estimate the
81 satellite remote sensing-derived sea surface $p\text{CO}_2$ data in the East China Sea during 2000–2014, and then used this algorithm to
82 estimate sea surface $p\text{CO}_2$ for the whole China Sea. These authors also pointed the limitation of their mechanistic
83 semi-analytical algorithm (MeSAA) which did not fully account for some local processes and therefore causes errors (the RMSE
84 is about 45 μatm in the SCS (reported in Wang et al., 2021). Bai et al. (unpublished) subsequently used a machine learning based
85 non-linear regression to develop a retrieval algorithm for seawater $p\text{CO}_2$ in the China Sea, and the satellite-derived $p\text{CO}_2$ data
86 from 2003-2018 were provided by the SatCO₂ platform (www.SatCO2.com) with a RMSE of 21.1 μatm in China Seas.

87 To take advantages of both the high spatiotemporal resolution of the remote sensing-derived $p\text{CO}_2$ data (RSdata) and the accuracy
88 of the observational data, Wang et al. (2021) reconstructed the basin-scale sea surface $p\text{CO}_2$ in the SCS in summer by using the
89 empirical orthogonal function (EOF) based on a multi-linear regression method and demonstrated the reliability of the
90 reconstructions. However, when the spatial standard deviation of observed data is relatively large because of the influence of
91 outliers, the reconstruction results may be biased (Wang et al., 2021). Therefore, many studies used the machine learning-based



92 regression method to reduce the influence of outliers for open ocean areas, with a RMSE of $<17 \mu\text{atm}$ in most cases (e.g., Zeng et
93 al., 2017; Li et al., 2019).

94 Building upon the EOF method that significantly improved the reconstruction in terms of spatial pattern and accuracy (Wang et al.,
95 2021), we developed a machine learning-based regression method facilitated by the EOF to fully resolve the long-term spatial
96 distribution of sea surface $p\text{CO}_2$ at a resolution of $0.05^\circ \times 0.05^\circ$ in the SCS. In addition to the typical machine learning performance
97 metrics, we present a novel uncertainty calculation method that incorporates the bias of both the reconstruction and the sensitivity
98 of reconstructed models.

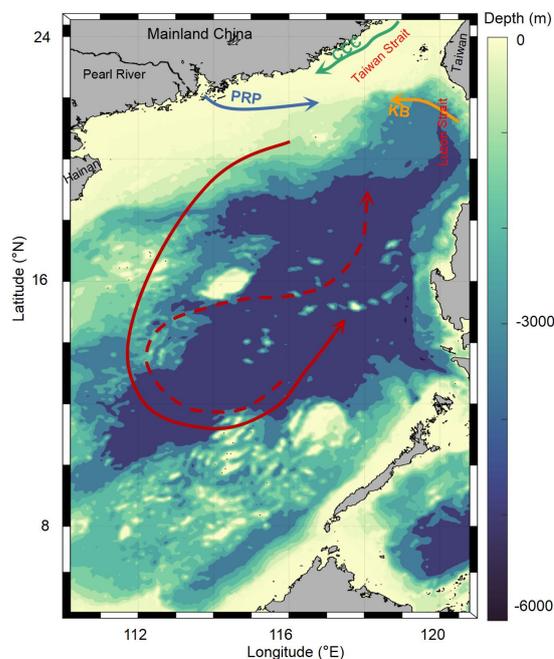
99

100 **2 Study site and data sources**

101 **2.1 Study area**

102 The SCS, located in the western Pacific, has a maximum water depth of ca. 4700 m (e.g., Gan et al., 2006, 2010). The rhombus
103 deep-water basin with a southwest-northeast direction accounts for about half of the total area of the SCS (Fig. 1). The
104 oceanography of the SCS is largely modulated by the Asian monsoon and the topography, thus exhibiting seasonally varying
105 surface circulation, river inputs, and upwelling. Forced by the northeast winds in winter, the circulation of the upper layer shows a
106 large cyclonic circulation structure (red solid line in Fig. 1), while in summer it exhibits an anticyclonic circulation structure
107 forced by southwest winds (red dashed line in Fig. 1; Hu et al. 2010). In the northern SCS, the Pearl River discharges into the SCS
108 with an annual freshwater input of $3.26 \times 10^{11} \text{ m}^3$ (e.g., Dong et al., 2004; Dai et al., 2014). The area influenced by the Pearl River
109 Plume may extend southeastward to a few hundred kilometers from the river estuary in summer because of the monsoon wind
110 stress (Dai et al., 2014). The northern and western coastal regions of the SCS also feature summer coastal upwelling in summer,
111 mainly including the Eastern Guangdong and Qiongdong upwelling systems in the northern SCS and the Vietnam upwelling
112 systems in the western SCS (e.g., Cao et al., 2011; Chen et al., 2012; Gan et al., 2006; Gan et al., 2010; Li et al., 2020).

113



114

115 **Figure 1. Topographic map of the South China Sea (SCS) showing a basin wide cyclonic gyre in winter (solid line) and an**
116 **anticyclonic gyre over the southern half of the SCS in summer (dashed line). Also shown are the Kuroshio Branch (KB,**
117 **orange line), the China Coastal Current (CCC, green line), and the Pearl River plume (PRP, blue line).**

118 The SCS is a semi-enclosed sea basin with dynamic water exchange with the East China Sea via the Taiwan Strait and Western
119 Pacific via the Luzon Strait (Fig. 1). In winter, driven by the winter monsoon, the China Coastal Current (CCC, yellow line in Fig.
120 1; Han et al., 2013; Yang et al., 2022) flows south along the Chinese mainland through the Taiwan Strait, and occupies the
121 northern SCS with cold, fresh, nutrient-rich waters. The strong northeast winds in winter also slow down the western boundary
122 ocean current, forcing the intrusion of Kuroshio water, which shows high surface salinity and high total alkalinity, into the SCS
123 via the Luzon Strait (orange line in Fig. 1; Du et al., 2013; Park, 2013; Yang et al., 2022).

124 **2.2 Observational $p\text{CO}_2$ data**

125 Data collected from field surveys during the study period 2003-2020 are summarized in Table 1. Most observations were made in
126 July, and fewer observations were made in March and December of each year. The rough sea state in the SCS in winter and early
127 spring limited the survey during these seasons. Data collected from July 2000 to January 2018 were originally published by Li et
128 al. (2020). The in situ $p\text{CO}_2$ were collected from R/Vs *Dongfanghong-2*, *Tan Kah Kee (TKK)*, etc. (shown in Table 1). The data
129 collection methods used in this study have been introduced in Li et al. (2020). The spatial coverage and frequency of the
130 observations are shown in Figure 2 and show that there are pronounced seasonal changes and that the data cover a large spatial
131 area. For example, the spatial coverage of the observed data in spring and fall are relatively uniformly distributed, and the south
132 end of the spatial coverage reaches 5°N in spring, whereas that during other seasons is concentrated in the northern and central

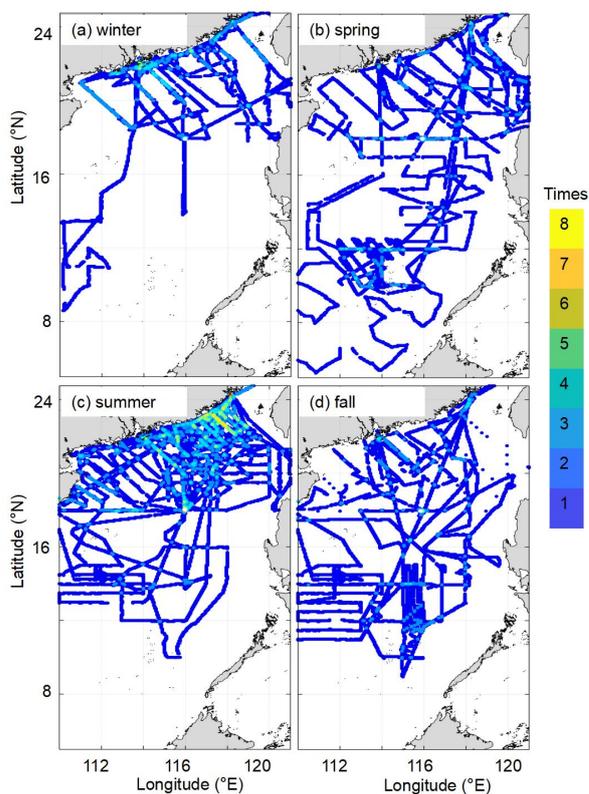


133 regions of the SCS. In addition, only one observation was made in the basin area in winter, while the northern coastal area was
 134 more frequently surveyed, especially in summer.

135 **Table 1. Summary of the seasonal observational data of sea surface $p\text{CO}_2$ in the South China Sea for the period 2003-2020**
 136 **used in this study.**

Season	Spring			Summer		
	March	April	May	June	July	August
Cruise time					2004.07	
		2005.04		2006.06	2005.07	
		2008.04	2004.05	2016.06	2007.07	2007.08
	2004.03	2009.04	2011.05	2017.06	2008.07	2008.08
		2012.04	2014.05	2019.06	2009.07	2019.08
		2020.04	2020.05	2020.06	2012.07	
					2015.07	
					2019.07	
Season	Fall			Winter		
	September	October	November	December	January	February
Cruise time	2004.09				2009.01	
	2007.09	2003.10	2006.11	2006.12	2010.01	2004.02
	2008.09	2006.10	2010.11		2018.01	2006.02
	2020.09					

137 ***Data were collected before February 2018 except those collected in July 2015 and June 2017 which are from Li et al.**
 138 **(2020).**



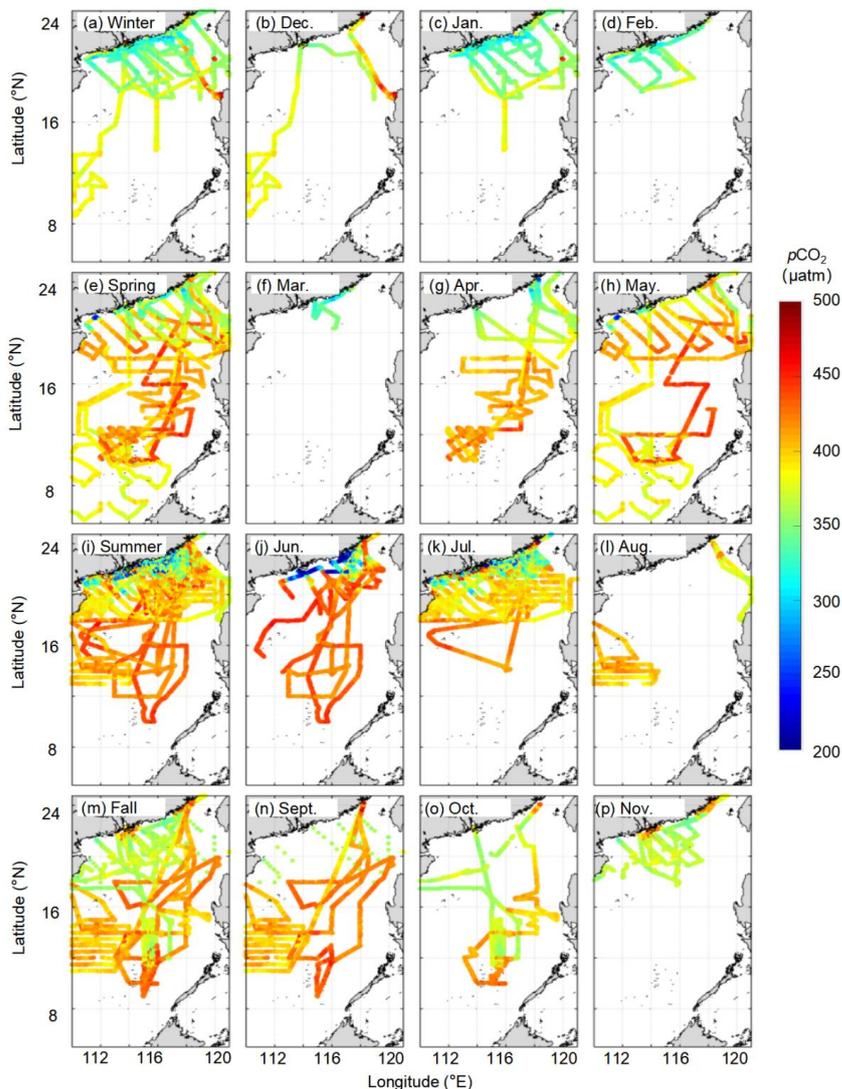
139

140 **Figure 2. Cruise tracks of the observations conducted in the South China Sea in each season from 2000 to 2020: (a) winter,**
141 **(b) spring, (c) summer, and (d) fall. The data were collected before February 2018 except those collected in July 2015, June**
142 **2017 which are from Li et al. (2020).**

143

144 Figure 3 shows the spatial and temporal distributions of surface water $p\text{CO}_2$. Seasonally, the lowest $p\text{CO}_2$ occurs in January, and
145 the highest concentrations occur in May and June. Spatially, the $p\text{CO}_2$ distribution in the basin is relatively homogeneous, but
146 shows large variability in the northern region. In the northern coastal area in summer, the observed $p\text{CO}_2$ distribution is affected
147 by the Pearl River plume (yielding low values) and coastal upwelling (yielding high values), which last into early fall. In winter
148 and early spring, relatively low $p\text{CO}_2$ values were determined in the coastal area. In addition, the high $p\text{CO}_2$ values recorded on
149 the western side of the Luzon Strait in December demonstrate the influence of winter upwelling during some of the surveys.

150 In addition to the above observational data, we selected four independent surveys corresponding to four seasons, and a continuous
151 observation dataset during 2003–2019 at the Southeast Asia Time-Series (SEATs) station (Dai et al., 2022) as two important
152 independent testing datasets.



153

154 **Figure 3. Seasonal and monthly sea surface $p\text{CO}_2$ fields in the South China Sea. The data sources can be found in Table 1.**

155

156 2.3 Remote sensing-derived sea surface $p\text{CO}_2$ data

157 The gridded ($0.05^\circ \times 0.05^\circ$) remote sensing-derived $p\text{CO}_2$ data covered almost the entire SCS ($5\text{--}25^\circ \text{N}$, $109\text{--}122^\circ \text{E}$), and show the
158 major CO_2 variation at a large scale (Wang et al., 2021; Bai et al., unpublished). Further details of the remote sensing (RS) data
159 can be found in the Sat CO_2 platform (www.SatCO2.com).

160 A grid-to-grid comparison was undertaken (Fig. 4) and the RMSE of the RS data-derived $p\text{CO}_2$ values were compared with the
161 observed $p\text{CO}_2$ data (Table 2). This comparison shows that the difference between the RS and the observed $p\text{CO}_2$ data ranges from



162 35 to 120 μatm in the coastal area, and that the largest biases occur in summer. In terms of the RMSE (Table 2), the largest bias
 163 reaches 30.0 μatm in summer. Bai et al. (2015; unpublished) and Wang et al. (2021) pointed out that relatively large discrepancies
 164 may reflect the limitations of the current algorithm, which considers only biological processes and the turbidity induced by the
 165 Pearl River discharge (characterized by Chl *a* and the remote sensing reflectance at 555 nm (*rrs*555) and does not take into
 166 account the riverine dissolved inorganic carbon and the input of other substances that may affect $p\text{CO}_2$.
 167 Because of the relatively large bias of RS data, their direct use may lead to overestimate $p\text{CO}_2$ values in the reconstructed field. In
 168 this study, the EOF method was used to compute the spatial patterns of the RS data to assess the accuracy and spatial distribution
 169 of the reconstruction data as a whole and remove the influence of the limited RS data.

170

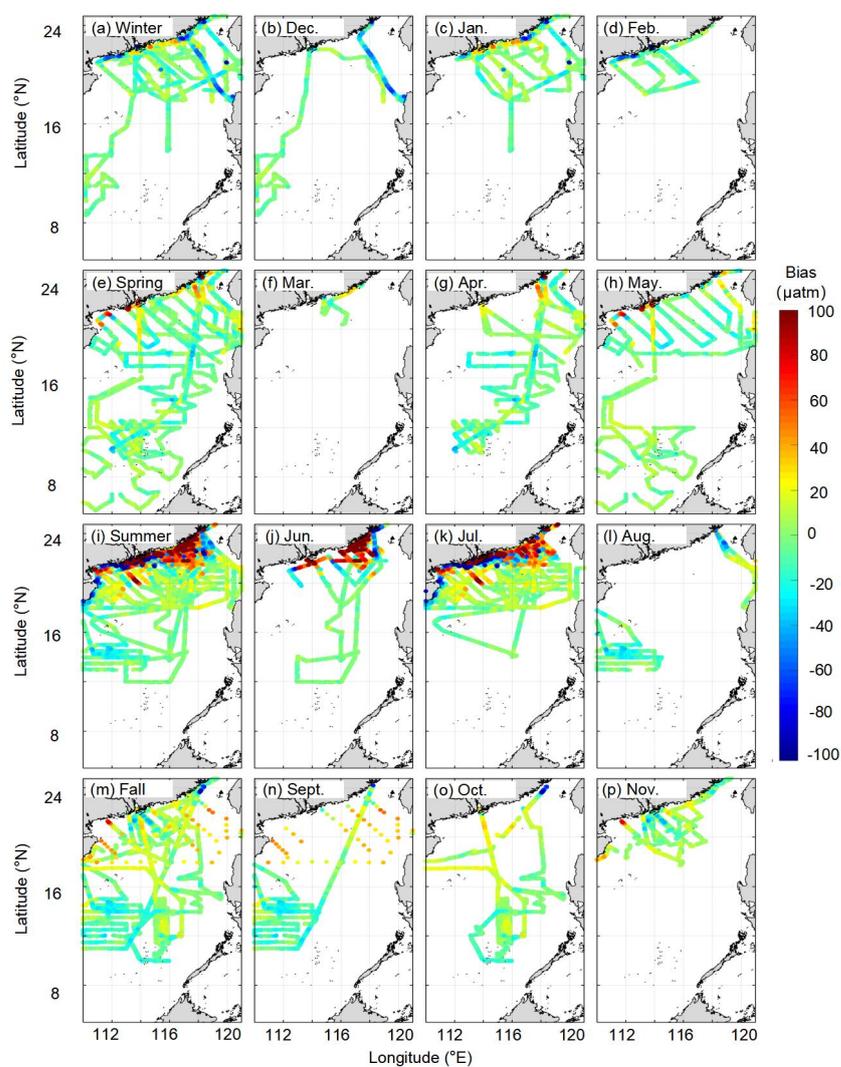
171 **Table 2. Biases between the seasonal remote sensing (RS)-derived $p\text{CO}_2$ data and observed $p\text{CO}_2$ data, and between the**
 172 **reconstructed and the observed underway $p\text{CO}_2$ data. (unit: μatm ; the remote sensing-derived $p\text{CO}_2$ data during**
 173 **2003-2019 are from www.SatCO2.com and the source of observed data can be found in Table1. The reconstructed $p\text{CO}_2$**
 174 **data are from section 3; all data were gridded into $0.05^\circ \times 0.05^\circ$; /: no data). MAE = mean absolute error; RMSE = root**
 175 **mean square error; R^2 = coefficient of determination; MAPE = mean absolute percentage error.**

	RS	Training data	Testing data I	Testing data II	Testing data III	
Spring	MAE	9.00	2.44	4.76	1.68	/
	RMSE	12.70	3.47	7.43	2.26	/
	R^2	/	0.98	0.92	/	/
	MAPE	/	0.01	0.01	/	/
Summer	MAE	16.75	2.48	8.46	5.73	/
	RMSE	29.956	3.54	14.69	15.18	/
	R^2	/	0.99	0.89	/	/
	MAPE	/	0.01	0.02	/	/
Fall	MAE	9.93	2.41	4.90	7.133	/
	RMSE	13.08	3.39	6.85	8.94	/
	R^2	/	0.98	0.92	/	/
	MAPE	/	0.01	0.01	/	/
Winter	MAE	9.25	2.18	5.61	11.41	/
	RMSE	14.26	3.14	8.82	12.63	/
	R^2	/	0.98	0.89	/	/



	MAPE	/	0.01	0.01	/	/
	MAE	11.95	2.41	6.30	5.27	6.19
Annual	RMSE	20.66	3.43	10.79	11.18	8.26
	R ²	/	0.99	0.91	/	/
	MAPE	/	0.01	0.01	/	/

176



177

178 **Figure 4.** Differences between the seasonal and monthly remote sensing-derived $p\text{CO}_2$ and the observed $p\text{CO}_2$ data (the
 179 former during 2003-2019 is from www.SatCO2.com, and the source of observed data can be found in Table 1. Both
 180 datasets are gridded into $0.05^\circ \times 0.05^\circ$, and the bias is plotted grid by grid).



181

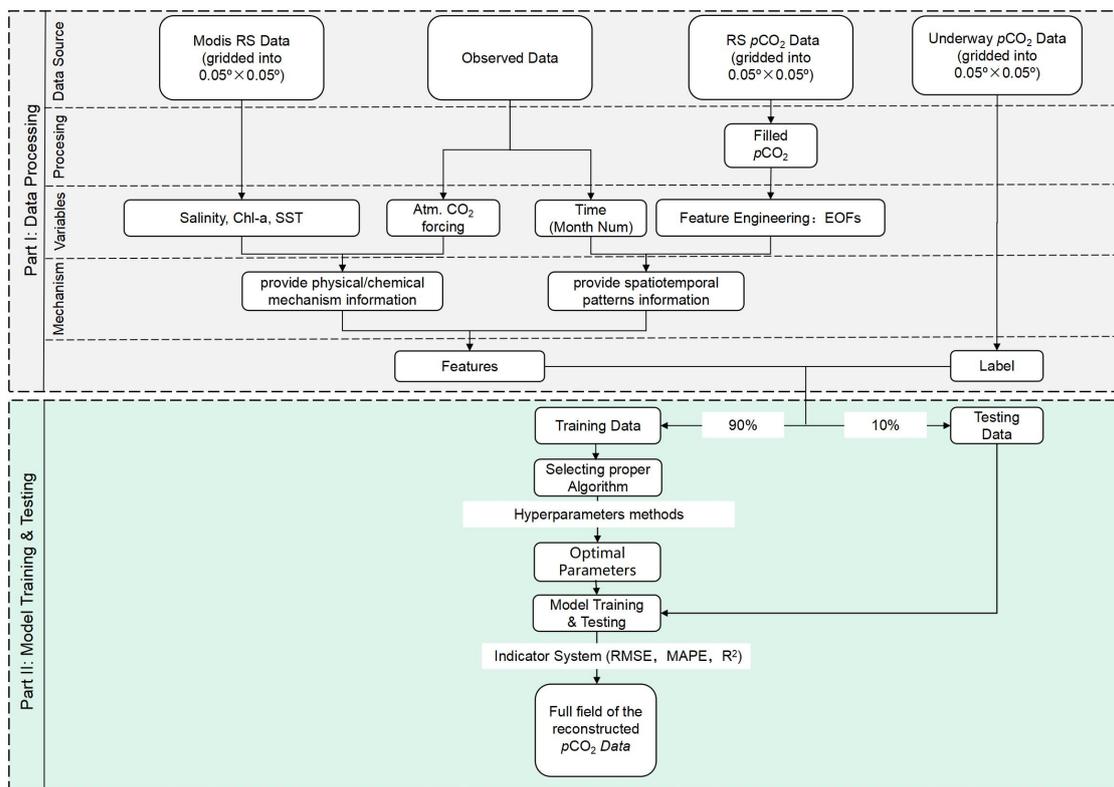
182 **2.4 Other data**

183 The RS SST data produced by MODIS (<https://oceancolor.gsfc.nasa.gov/>) are adopted here. The uncertainty of this dataset in the
184 SCS is $\sim 0.27^\circ$ (Qin et al., 2014). For the SSS data, Wang et al. (in preparation) found a relatively high differential between the
185 values found in different open source databases (i.e., multi-satellite fusion data from <https://podaac.jpl.nasa.gov/>; model data from
186 <https://climatedataguide.ucar.edu/>; multidimensional covariance model data from <https://resources.marine.copernicus.eu/>) and our
187 observed data. Wang et al. (in preparation) reconstructed a SSS database using machine learning methods based on the observation
188 dataset. The bias between the reconstructed SSS and observed data was near-zero (mean absolute error, MAE: ~ 0.25). Next, we
189 used Chl-*a* as an indicator of biological influence. Chl-*a* data from MODIS (<https://oceancolor.gsfc.nasa.gov/>) are adopted in the
190 present study, which have a bias of ~ 0.35 on log scale and $\sim 115\%$ in the SCS (Zhang et al., 2006). Atmospheric $p\text{CO}_2$ also
191 influences surface water $p\text{CO}_2$ through air–sea CO_2 exchange. We chose the atmosphere CO_2 mole fraction ($x\text{CO}_2$) data from the
192 monthly mean CO_2 concentrations measured at Mauna Loa Observatory, Hawaii (<https://gml.noaa.gov/>). The atmospheric $p\text{CO}_2$
193 values were calculated from $x\text{CO}_2$ using the method of Li et al. (2020).

194

195 **3 Methods**

196 The $p\text{CO}_2$ reconstruction procedure is shown in Figure 5. It includes: (1) data processing and (2) model training and testing. For
197 the former, we first gridded the observed and RS $p\text{CO}_2$ data into $0.05^\circ \times 0.05^\circ$ grid boxes with monthly temporal resolution.
198 Secondly, we used the $p\text{CO}_2$ filling method of Fay et al. (2021) to fill the missing $p\text{CO}_2$ measurements with the RS $p\text{CO}_2$ data. We
199 then used a feature engineering method to ignore any biases in the data itself or from the $p\text{CO}_2$ filling method. Thirdly, the gridded
200 observed $p\text{CO}_2$ data and their corresponding RS data were divided into a training set (90%) and a testing set (10%) to calculate the
201 $p\text{CO}_2$ reconstruction model. To ensure that the model had sufficient training samples in the coastal area, we divided the entire SCS
202 into two regions along the 200 m depth contour. The data from these two regions were divided into training and testing sets with
203 the same ratios listed above (9:1), which were then combined to obtain the final training and testing sets.



204

205 **Figure 5. Procedure for the reconstruction of surface water pCO_2 using machine learning. RS data = remote sensing data,**
 206 **RMSE = root mean square error, MAPE= mean absolute percentage error, and R^2 = coefficient of determination, and**
 207 **MAE = average absolute error.**

208 For model training and testing, we first chose a relatively reliable algorithm to undertake the pCO_2 reconstruction. After that, we
 209 determined the optimal range of the parameters using hyperparameter methods (code from <https://github.com/optuna/>) for the
 210 training set. The final optimal parameter values were then determined using the K-fold and cross validation method (code from
 211 <https://github.com/suryanktiwari/Linear-Regression-and-K-fold-cross-validation>) for the training set. These optimal parameters
 212 were applied to the chosen algorithm. Finally, the testing set was used to verify the accuracy of the pCO_2 reconstruction model
 213 produced by the training set, and some indicators of the model's accuracy were calculated. More detailed methods employed in
 214 the present study are described below.

215 3.1 Remote sensing data filling

216 As mentioned in the SatCO2 platform (www.SatCO2.com), the RS pCO_2 data are missing some values. Thus, we used the pCO_2
 217 filling method suggested by Fay et al. (2021) to fill in the missing portions. First, a scaling factor for a filled month was calculated
 218 according to Equation 1:

$$219 \quad sf_{pCO_2} = \text{mean}_{x,y} \left(\frac{pCO_2^{ens}}{pCO_2^{clim}} \right) \quad (1)$$



220 where sf_{pCO_2} is the scaling factor, pCO_2^{ens} is the monthly RS pCO_2 datum, and pCO_2^{lim} is the monthly climatology RS pCO_2
221 datum; x and y indicate that we took the area-weighted average over longitude (x) and latitude (y) to produce the monthly sf_{pCO_2}
222 value. Then, the filled portion of the data can be calculated from the pCO_2^{lim} data multiplied by the sf_{pCO_2} value (see Fay et al.
223 (2021) for details of this method).

224 Briefly, this filling method scales the climatological monthly pCO_2 field values to fill in the missing measurements. Therefore,
225 although specific values may be biased, the interpolated measurements still retain the main spatial distribution pattern of the filled
226 months.

227 3.2 Feature engineering and selection

228 As mentioned above, the pCO_2 filling method may bias some of the actual values. To avoid such biases, instead of directly using
229 the RS pCO_2 data as features in our reconstructed model, we used the feengineered featured data (via the EOF method) to obtain
230 the main spatiotemporal distribution patterns of the RS pCO_2 data as features in our reconstructed model. The EOF method can
231 reflect the spatial commonality of variables shown in the time series. For each 12 months, the cumulative variance contribution of
232 the first eight EOF values was consistently $> 90\%$, indicating it that it could explain the main pCO_2 spatial characteristics during
233 each month, and we therefore selected them as features.

234 The feature selection in our reconstructed model can be divided into two main categories. The first one is related to the underlying
235 physicochemical mechanism controlling the pCO_2 distribution, and the other one can provide spatiotemporal information for
236 pCO_2 reconstruction. For example, the SST dominating the seasonal variation in surface water pCO_2 in the northern SCS (Zhai et
237 al., 2005; Chen et al., 2007; Li et al., 2020). Previous research (Landschützer et al., 2014; Laruelle et al., 2017; Denvil et al., 2019)
238 shows that Chl- a plays a critical role in fitting the influence of biological activity to pCO_2 , especially in the northern SCS
239 (Landschützer et al., 2014; Laruelle et al., 2017; Denvil et al., 2019). Sutton et al. (2017) suggest that the increase in atmospheric
240 pCO_2 controls the increase in seawater pCO_2 . For the features that provide spatiotemporal information for pCO_2 reconstruction,
241 whereas in the present study we selected the first eight EOF values of pCO_2 as the main spatial distribution feature and monthly
242 information of the observed datasets as the temporal feature.

243 3.3 Algorithm selection

244 Ensemble learning provides one of the most powerful machine learning techniques (e.g., Zhan et al., 2022; Chen et a., 2020). It is
245 the process of training multiple machine learning models and combining their output to improve the reliability and accuracy of
246 predictions (e.g., Zhan et al., 2022; Chen et a., 2020). Different models are used as the basis to develop an optimal predictive
247 model. There are two main ways to employ ensemble learning: bagging (to decrease the model's variance) or boosting (to
248 decrease the model's bias). The random forest algorithm (code from <https://scikit-learn.org/stable/>) is an extension of the bagging
249 method as it utilizes both bagging and feature randomness to create an uncorrelated forest of decision trees. The Light Gradient
250 Boosting Machine (LightGBM; code from <https://github.com/microsoft/LightGBM/>) is a gradient boosting framework that uses



251 tree-based learning algorithms. LightGBM can be used for regression, classification, and other machine learning tasks; it exhibits
252 rapid and high-performance as a machine learning algorithm. CATBOOST (code from <https://github.com/catboost/>) is a gradient
253 boosting algorithm, which improves prediction accuracy by adjusting weights according to the data distribution and by
254 incorporating prior knowledge about the dataset. This can help to reduce overfitting and improve generalization performance.
255 From the above options, we chose three ensemble learning algorithms as the machine learning-based regression portion, and
256 multi-linear regression methods (Wang et al., 2021) as the linear regression portion, and we then used the K-fold and cross
257 validation methods to verify the applicability of the different regression algorithms in the $p\text{CO}_2$ reconstruction of summer training
258 data in the SCS, since the greatest temporal sampling coverage occurs in summer (Table 1; Fig. 2). Results show that the
259 CATBOOST algorithm yields the best degree of accuracy, with an RMSE of 16 μatm ; for comparison, the RMSE of LightGBM
260 was 27 μatm , that of Random Forest was 26 μatm , and nearly 20 μatm was found for the linear regression algorithm employed by
261 Wang et al. (2021). Thus, CATBOOST appears to provide a relatively reliable algorithm for $p\text{CO}_2$ reconstruction.

262

263 3.4 Evaluation metrics

264 It is necessary to evaluate the accuracy of any model based on certain error metrics before applying it to specific scenarios.
265 Common model evaluation metrics include RMSE, MAPE, R^2 (coefficient of determination), and MAE.

266 The mean squared error (MSE) stands for the standard deviation of the residuals (prediction error), where the residuals represent
267 the distance between the fitted line and the data point i.e., it stands for the degree of concentration of the reconstructed data around
268 the regression line. In regression analysis, RMSE is commonly used to verify experimental results. To assess bias, the RMSE
269 needs to combine the magnitude of the model data and is calculated as:

$$270 \quad \text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - y_{ri})^2} \quad . \quad (2)$$

271 where y stands for the observational data, y_r represents the reconstructed data, and n is the number of data.

272 The mean absolute percentage error (MAPE) is a statistical measure used to define the accuracy of a machine learning algorithm
273 on a particular dataset. It is commonly used because, compared to other metrics, it uses a percentage to measure the magnitude of
274 the bias and is easy to understand and interpret; the lower the value of MAPE, the better a model is at forecasting. MAPE is
275 calculated as follows:

$$276 \quad \text{MAPE} = \frac{1}{n} \sum_{i=1}^n \frac{|y_i - y_{ri}|}{|y_i|} \quad (3)$$

277 The regression error metric, the coefficient of determination R^2 , can describe the performance of a model by evaluating the
278 accuracy and efficiency of modeled results, i.e., it indicates the magnitude of the dependent variable calculated by the regression
279 model that can be explained by the independent variable, and is calculated as:



280
$$R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{\sum_{i=1}^n (y_i - y_{ri})^2} \quad (4)$$

281 MAE is the average absolute difference between the field observations (true values) and model output (predicted values). The sign
282 of these differences is ignored so that cancellations between positive and negative values do not occur. It is calculated as:

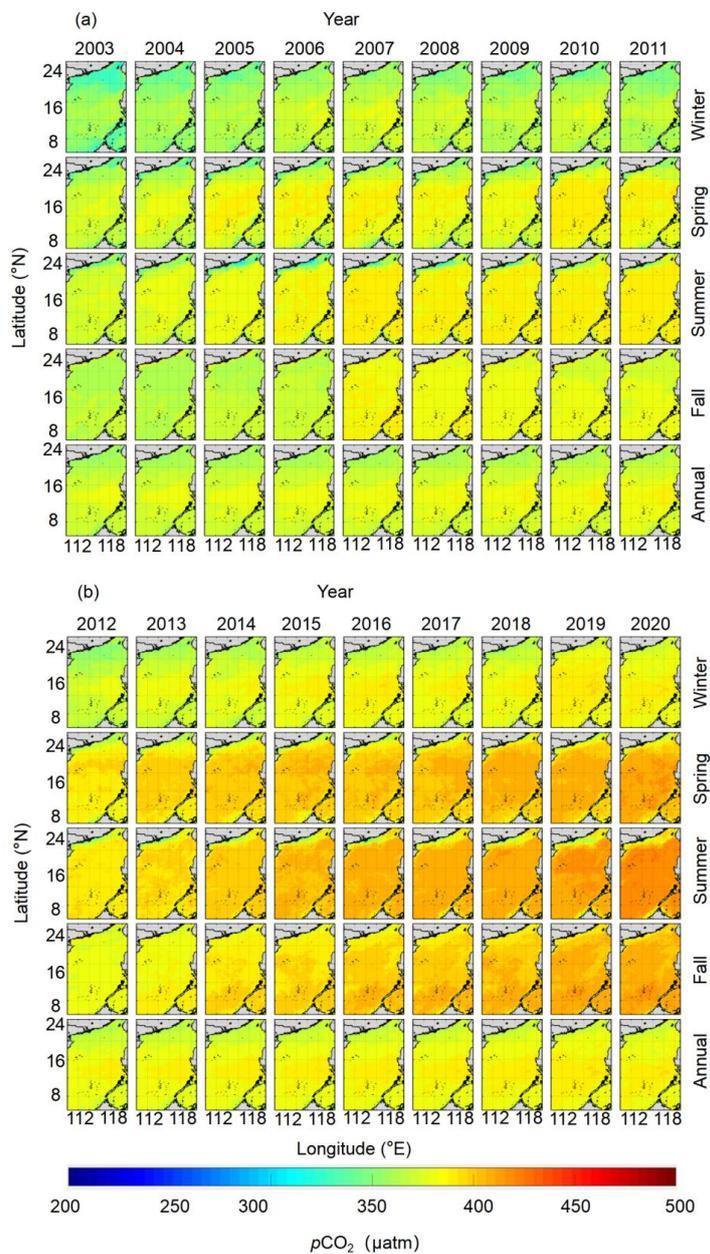
283
$$MAE = \frac{1}{n} \sum_i^n |y_i - y_{ri}| \quad (5)$$

284

285 **4 Results and discussion**

286 **4.1 Results**

287 The reconstructed $p\text{CO}_2$ fields show relatively low values in the northern coastal study region but generally shows high values in
288 the mid and southern basins (Fig. 6). The continuity changes of the spatiotemporal distribution can be found in the reconstruction
289 results (Fig. 6). The reconstructed $p\text{CO}_2$ fields show a trend of slow but sustained increase from 2003 to 2020. Spatial patterns of
290 $p\text{CO}_2$ change between 2003 and 2020, such that the coastal portion of the northern SCS shows relatively complex variability
291 because of multiple controlling factors, such as coastal upwelling, river plumes, biological activity, etc. However, $p\text{CO}_2$ values in
292 the mid and southern basin are relatively homogeneous, because they are mainly controlled by atmospheric CO_2 forcing and SST.
293 Temporal changes in $p\text{CO}_2$ between 2003 and 2020, are relatively large ($\sim 44 \mu\text{atm}$) in summer and relatively small ($\sim 33 \mu\text{atm}$) in
294 winter.



295

296 **Figure 6. Reconstructed seasonal and annual $p\text{CO}_2$ fields in the South China Sea during the period 2003 to 2020 (a,**
297 **2003-2011; b, 2012-2020).**

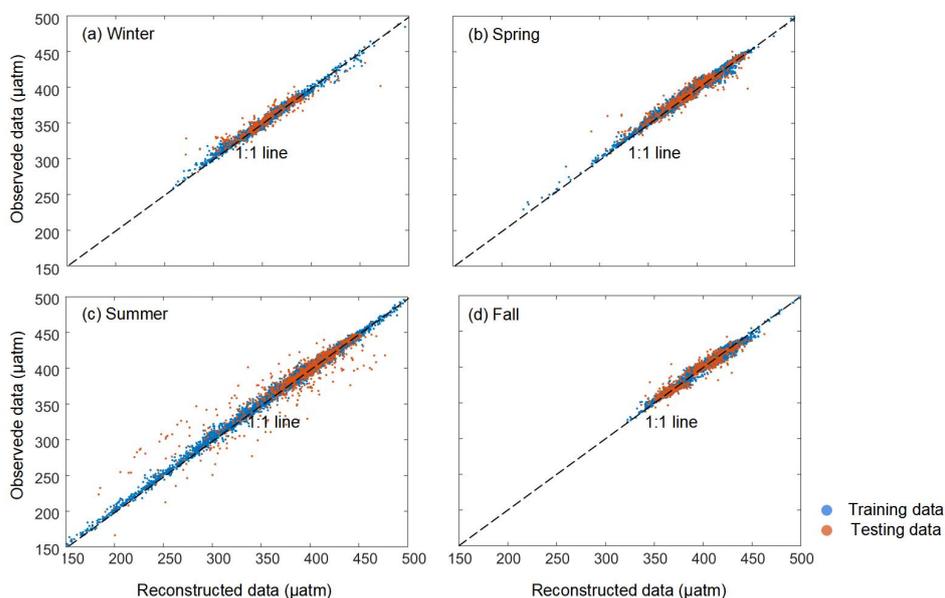
298 4.2 Model validation

299 Figure 7 compares the reconstructed and field-observed data. For the training dataset, the reconstructed $p\text{CO}_2$ fields of the four
300 seasons fit the field-observed data well (Fig. 7), with an average RMSE of 3.43 μatm and an average MAE of ~ 2 μatm (Table 2).

301 For the testing sets, although there are some outliers, most of the reconstructed $p\text{CO}_2$ data are consistent with field-observed data,



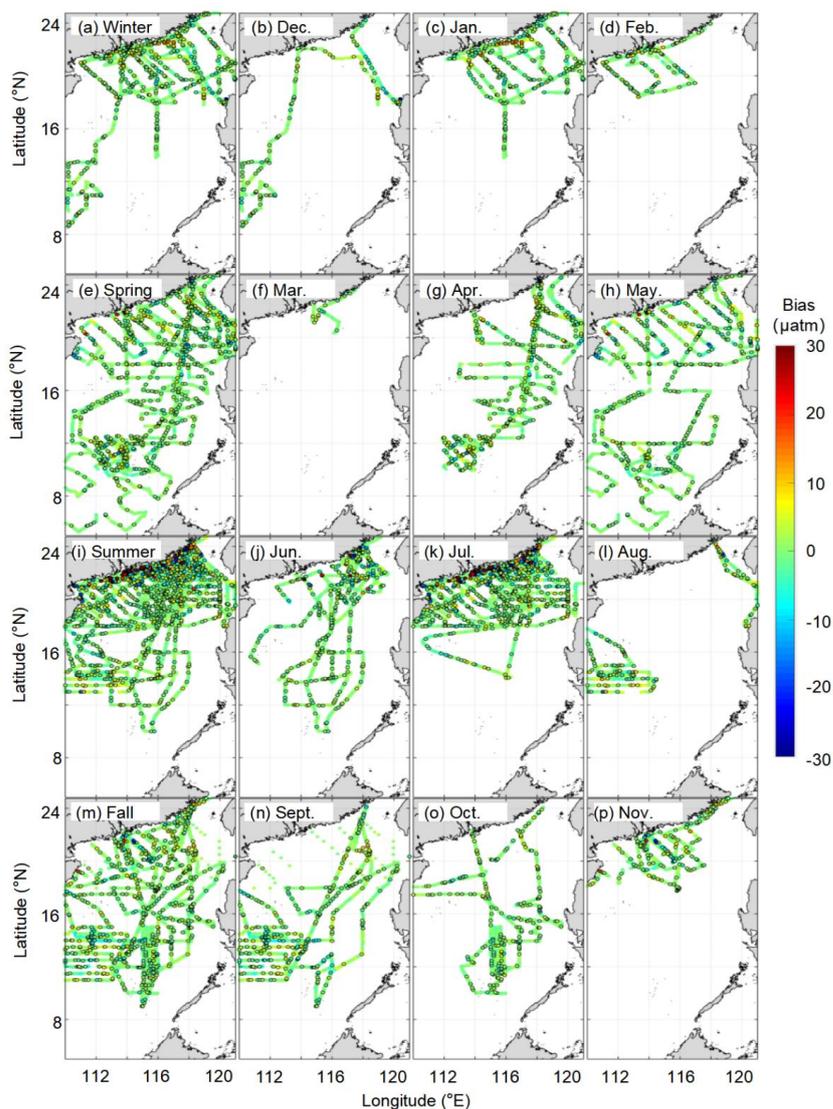
302 with RMSE averaging $10.79 \mu\text{atm}$ and MAE averaging $6.30 \mu\text{atm}$. The R^2 of the testing set is ca. 0.91. In terms of MAPE, the
303 accuracies of the four seasonal models are all around 99% (Table 2), with the highest value for spring data and the lowest value
304 for summer data. The greatest bias in the summer may be the influence of relatively complex regional processes, such as river
305 plumes and upwelling. The four evaluation metrics indicate that our reconstructed $p\text{CO}_2$ field is highly accurate in simulating both
306 the training and testing sets.



307
308 **Figure 7. Comparison between the reconstructed and the observed $p\text{CO}_2$ values.**

309 The distribution pattern of the biases between the reconstructed fields and the field observations in both training and testing
310 datasets can be found in Figure 8. In terms of the temporal distribution pattern, the biases are concentrated mainly in summer. For
311 the spatial distribution pattern, the biases in the northern coastal area are much greater than those in the basin. However, 95% of
312 the biases are $< \pm 10 \mu\text{atm}$. Therefore, our reconstruction data exhibit relatively high accuracy.

313



314

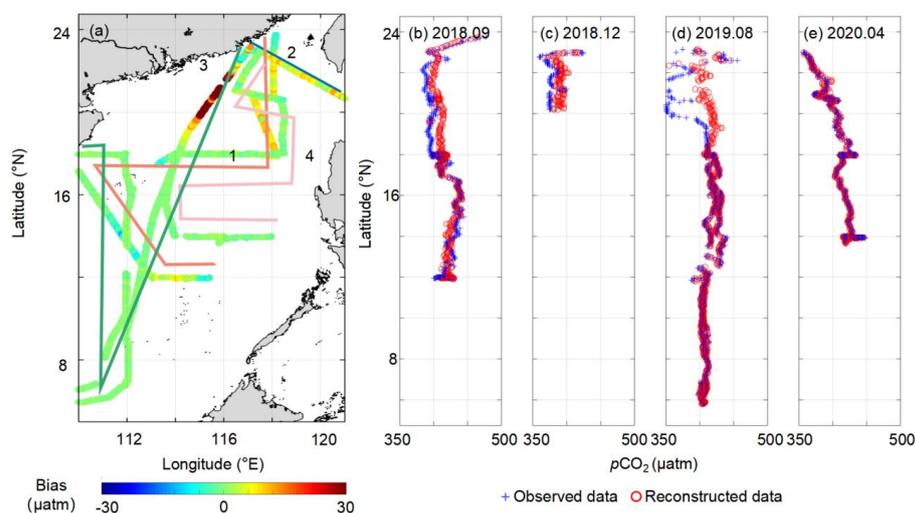
315 **Figure 8. Differences between the seasonal and monthly reconstructed $p\text{CO}_2$ and the observed $p\text{CO}_2$ data. The open circles**
316 **represent the difference between the training set and observational data, and the solid black circles represent difference**
317 **between the test set and observational data.**

318 .

319 Figure 9 shows the bias between our reconstructed fields and the four independent field observation datasets corresponding to the
320 four seasons. This validation can verify the accuracy of the reconstruction model in data months with no observations, namely the
321 applicability of the reconstruction model extrapolation. This comparison shows that the reconstruction model is relatively
322 accurate in the basin, with a near-zero bias (MAE: $\sim 8 \mu\text{atm}$, Fig. 9 a). The greatest bias occurs in the Pearl River plume area in



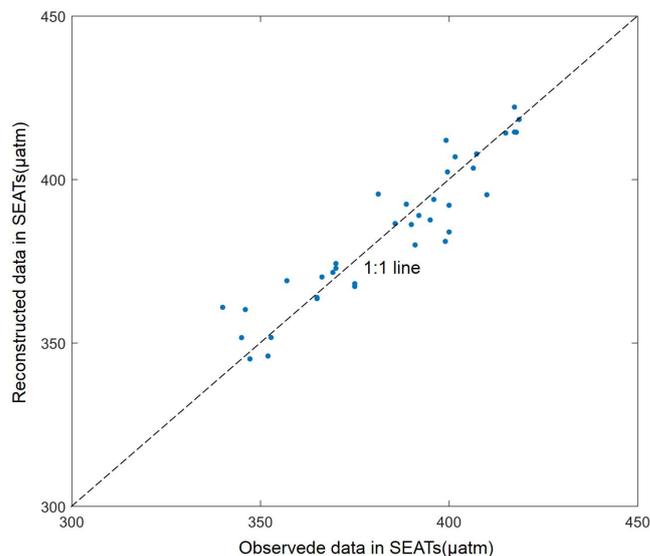
323 summer. The reconstruction model also has high accuracy in $p\text{CO}_2$ spatial variation trends, except in the Pearl River plume area
324 in summer (22–20 °N), as shown in Fig. 9 b–e). The effect of the Pearl River plume on the $p\text{CO}_2$ spatial distribution in our
325 reconstruction model is smaller than that shown by the field-observed data. This is because at around the survey time (August
326 24–28, 2019), a large amount of precipitation (~30mm/day; <https://psl.noaa.gov/data/gridded/data.ncep.reanalysis2.surface.html>)
327 occurred around the Pearl River estuary region (24–20 °N), which led to intensification of the Pearl River plume, such that the
328 plume with relative low CO_2 values eventually decreased the observed values. However, the monthly average runoff of the Pearl
329 River during that month (August, 2019; <http://www.pearlwater.gov.cn/>; Pearl River Plume Index in Wang et al., in preparation) is
330 low, indicating that our reconstruction model is still highly reliable from the a monthly average perspective. Thus, the
331 inconsistency between the reconstructed (monthly average) and field-observed datasets is mainly due to the differences in the time
332 scales of the remote sensing and the field-observed data. The reconstructed data in this study were determined on a monthly scale,
333 while the temporal resolution of the field-observed data was on the order of hours. It is clear that relatively pronounced short-term
334 changes in $p\text{CO}_2$, such as the diurnal variation caused by short-term heavy precipitation, cannot be reflected in the reconstructed
335 data.



336

337 **Figure 9. Difference between the reconstructed $p\text{CO}_2$ data and four independently observed datasets during the four**
338 **seasons. In (a), the numbers 1–4 represent September (2018.9), b), December 2018 (2018.12, c), August 2019 (2019.8, d), and**
339 **April 2020 (2020.4, e), respectively.**

340 Dai et al. (2022) produced a time-series of observed data from 2003 to 2019 at the SEATs station, which we used here to validate
341 the accuracy of the long-term trends of our model data (results shown in Fig. 10). The long-term trend of reconstructed $p\text{CO}_2$
342 data at the SEATs station are largely consistent with the observations, with differences mainly found before 2005. Thus, the
343 long-term trend of our reconstructed model is also highly reliable.



344

345 **Figure 10. Comparison of the reconstructed $p\text{CO}_2$ with the observations at the Southeast Asia Time Series (SEATs) station**
 346 **(116° E, 18° N). The observed data are from Dai et al. (2022), which were calculated from dissolved inorganic carbon and**
 347 **total alkalinity.**

348

349 4.3 Uncertainties

350 In previous studies, RMSE and MAE were mostly used to represent the uncertainties in the reconstructed data. As shown in Table
 351 2, our reconstruction data have a high degree of accuracy, with an RMSE of $\sim 10 \mu\text{atm}$ and MAE of $\sim 6 \mu\text{atm}$. However, this
 352 expression of uncertainty ignores the sensitivity of the reconstructed model to the features; i.e., the bias that the features
 353 themselves pass to the reconstructed model are ignored. Moreover, it is clearly unreasonable to use a single RMSE or MAE value
 354 to represent the entire region because the spatial bias patterns clearly differ between coastal and basin areas (Figs. 8 and 9).

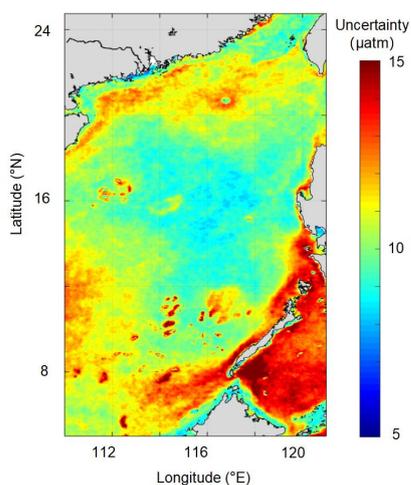
355 Thus, we here present a novel method of uncertainty calculation as shown below:

$$\begin{aligned}
 \text{356 } \text{Uncertainty} = \text{MAX} & \left(\frac{\sum_{i=1, j=1}^n \frac{|OR_Monthly_Data(i, j) - Obs_Monthly_Data(i, j)|}{Obs_Monthly_Data(i, j)}}{\text{num}(i) + \text{num}(j)} \right) * 100\% * p\text{CO}_2_recon \quad (\text{part 1}) \\
 \text{357 } & + \left(\frac{\partial p\text{CO}_2}{\partial \text{Feature}} \right) d\text{Feature} \quad (\text{part 2}). \quad (6)
 \end{aligned}$$

358 Equation (7) includes two parts; the first is the conservative bias between the reconstructed $p\text{CO}_2$ fields and the observations (part
 359 1), and the second is the sensitivity of the reconstructed model to the features (part 2). For part 1, $OR_Monthly_Data(i, j)$ stands
 360 for the monthly reconstructed data at longitude(i) and latitude(j), and $Obs_Monthly_Data(i, j)$ stands for the monthly observed
 361 data at longitude (i) and latitude (j). Therefore, the conservative bias is the maximum value of the monthly error ratio between the
 362 observed data and the reconstructed data. For part 2, where $d\text{Feature}$ stands for the bias of the features, we conducted a
 363 sensitivity analysis using a chain rule to evaluate the influence of these bias in the features on $p\text{CO}_2$. The bias of RS $p\text{CO}_2$ is ~ 21



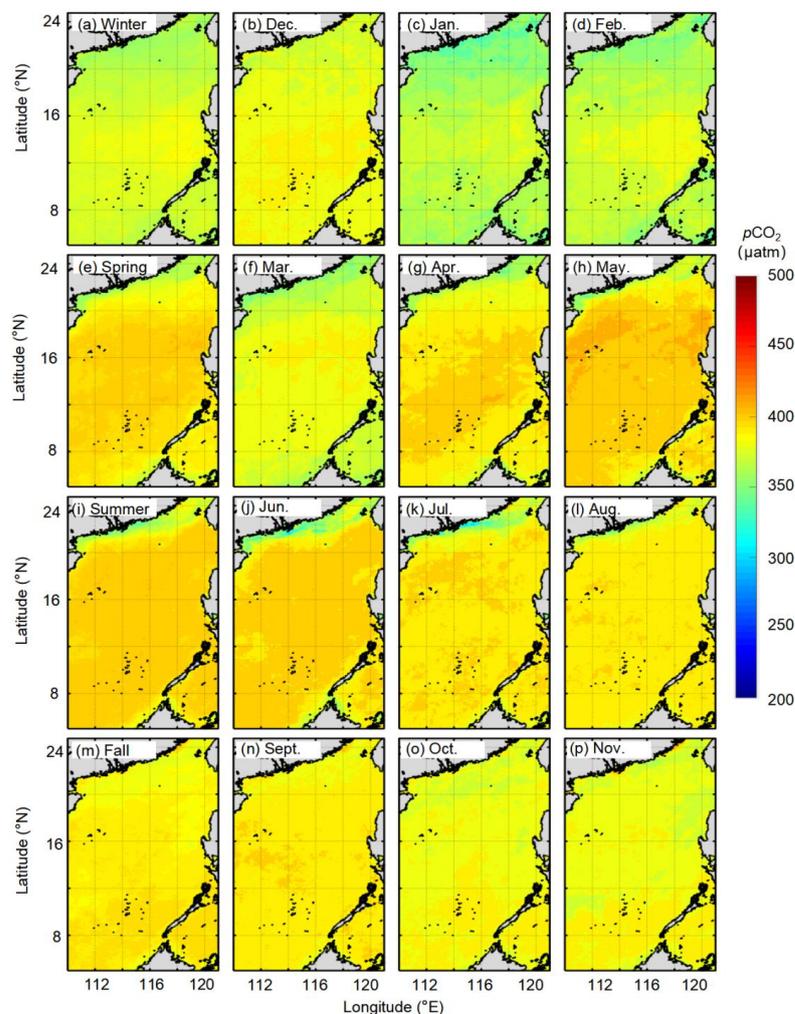
364 μatm (Table 2), the bias of SST is $\sim 0.27^\circ$ (Qin et al., 2014), that of SSS is ~ 0.33 (Wang et al., in prep.), and that of Chl-*a* is
365 $\sim 115\%$ (Zhang et al., 2006). We then estimated $p\text{CO}_2$ changes due to these features' variability by constraining these features
366 based on our model, and computing $\frac{\partial p\text{CO}_2}{\partial \text{Feature}}$. For example, for the $\frac{\partial p\text{CO}_2}{\partial \text{SST}}$ part, we only changed the value of SST, and kept the
367 value of the other features constant, to calculate the effect of each additional unit of SST on the results of the $p\text{CO}_2$ simulation.
368 These two parts were then added together to obtain the final uncertainty, and results are displayed in Figure 11. The uncertainties
369 are greater in the coastal area ($\sim 13 \mu\text{atm}$), than in the basin ($\sim 10 \mu\text{atm}$). The spatial pattern of the uncertainty is consistent with
370 that shown in Section 4.2.



371
372 **Figure 11. Uncertainties of the reconstructed $p\text{CO}_2$ fields.**

373
374 **4.4 Spatial and temporal $p\text{CO}_2$ features**

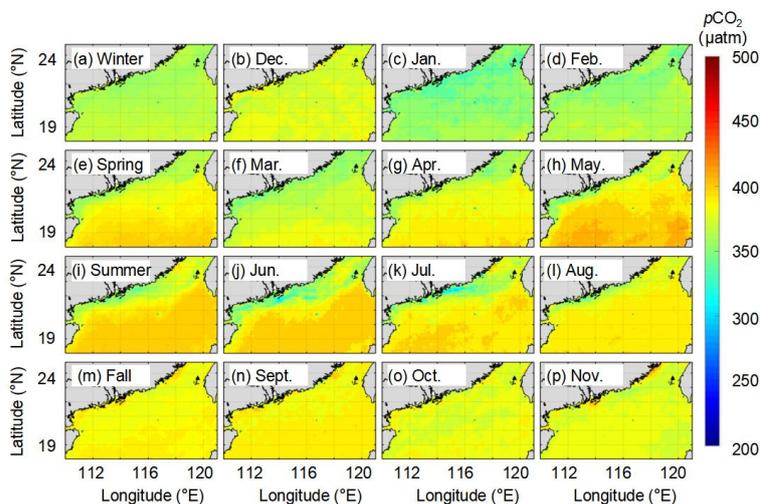
375 The climatological monthly reconstructed $p\text{CO}_2$ fields are shown in Figure 12. The highest values of the reconstructed $p\text{CO}_2$ fields
376 occur in May and June, and the lowest value occurs in January. In winter, $p\text{CO}_2$ first decreases in December and then increases in
377 January; the $p\text{CO}_2$ value is ca. $325 \mu\text{atm}$ in the northern coastal area, and ca. $350 \mu\text{atm}$ in the basin. In spring, $p\text{CO}_2$ gradually
378 increases from the basin to the northern coastal area, and the basin high-value center gradually expands outward starting in April.
379 In summer, $p\text{CO}_2$ gradually declines starting in June. In fall, $p\text{CO}_2$ increases from north to south, and the southern region shows
380 consistently high values.



381

382 **Figure 12. Long-term (2003–2020) seasonal and monthly average $p\text{CO}_2$ field (unit: μatm).**

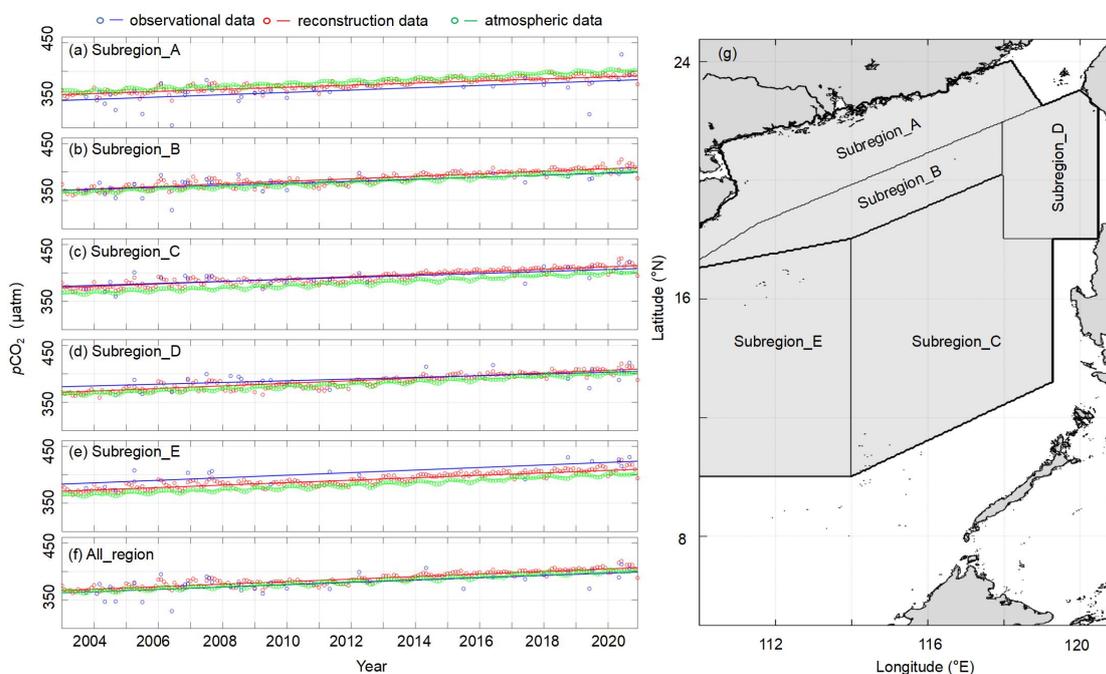
383 To better show specific regions in the northern coastal area, we zoomed in on the reconstructed $p\text{CO}_2$ fields at locations north of
384 18°N (Fig. 13). The reconstructed $p\text{CO}_2$ fields successfully reflect the influence of the meso-small scale processes on $p\text{CO}_2$ in this
385 northern coastal area of the SCS. For example, in winter, the relatively low $p\text{CO}_2$ values, which last into early spring, are mainly
386 controlled by the low SST, and the high $p\text{CO}_2$ around Luzon Strait affected by winter upwelling. In summer, the reconstructed
387 $p\text{CO}_2$ field shows that the influence of the Pearl River plume on $p\text{CO}_2$ is the strongest in July and lasts until September; it also
388 effectively shows the influence of coastal upwelling in the northeastern shelf ($\sim 23^\circ\text{N}$, 117°E). Thus, our reconstructed $p\text{CO}_2$ fields
389 clearly reflect the spatial pattern of the field observed $p\text{CO}_2$ (Fig. 3), which are generally consistent with previously reported
390 patterns (Li et al., 2020; Zhai et al., 2013; Gan et al., 2010).



391

392 **Figure 13. Long-term (2003–2020) seasonal and monthly averaged $p\text{CO}_2$ field in the region north of 18°N (unit: μatm).**

393



394

395 **Figure 14. Time series of Spatially averaged monthly $p\text{CO}_2$ data in five subregions (a-e) and the entire South China Sea (f)**
 396 **under study. The subregions are shown in (g). The lines indicate the deseasonalized long-term trend of the spatially**
 397 **averaged monthly $p\text{CO}_2$ data for each sub-region with the slopes shown in Table 3. The deseasonalized method can be**
 398 **found in Landschützer et al., 2016.**



399

400 **Table. 3 Deseasonalized long-term trend of the spatially averaged monthly $p\text{CO}_2$ data for each sub-region of the South**
401 **China Sea. (unit: $\mu\text{atm yr}^{-1}$).**

	All_region	A_region	B_region	C_region	D_region	E_region
Reconstructed	2.12±0.17	1.82±0.14	2.23±0.12	2.17±0.12	2.20±0.13	2.16±0.13
Observation	2.10±0.79	1.80±0.86	1.73±0.84	1.81±0.85	1.41±1.16	2.13±1.10

402

403 We divided SCS into five sub-regions according to Li et al. (2020). In Fig.14, region A stands for the northern coastal area of the
404 SCS, region B stands for the slope area of the northern SCS, region C stands for the SCS basin, region D stands for the region
405 West of the Luzon Strait, and region E stands for the slope and basin area of the western SCS. “All_region” indicates the whole
406 region containing the five sub-regions described above. We then calculated the deseasonalized long-term trend of spatially
407 averaged monthly data for each sub-region, and the results are shown in Figure 14 and Table.3. This deseasonalized trend is
408 consistent with that of observational data, and its uncertainty is on the 95% confidence interval much lower than that shown by the
409 observational data. We can thus also infer that the long-term trend of our reconstructed data shows high reliability in all
410 sub-regions, and that our data can serve as an important basis for predicting future changes of $p\text{CO}_2$ in the SCS.

411 In Fig.14 a-e, we found that the sea surface $p\text{CO}_2$ of the entire SCS is slightly higher than the atmospheric $p\text{CO}_2$, indicating that
412 the SCS is a weak source of atmospheric CO_2 . This conclusion is consistent with previous studies (e.g., Li et al., 2020). Moreover,
413 compared to the rate of atmospheric CO_2 increase ($\sim 2.2 \mu\text{atm yr}^{-1}$), for region A, the $p\text{CO}_2$ trend is much slower than that of
414 atmospheric $p\text{CO}_2$, and the spatially averaged monthly mean $p\text{CO}_2$ is lower than the atmospheric $p\text{CO}_2$. Thus, carbon
415 accumulation in this region is expected to increase in future. For regions C and E, the spatially averaged monthly mean $p\text{CO}_2$ is
416 higher than the atmospheric $p\text{CO}_2$; thus, these two regions will still provide a weak source of atmospheric CO_2 in future. Finally,
417 whether regions B and D act as a source or sink of atmospheric CO_2 is influenced by seasonal changes and physical processes.

418

419 **5 Data availability**

420 The data (the reconstructed CO_2 data, the Observational CO_2 data, and the remote sensing derived CO_2 data) for this paper are
421 available under the link <https://github.com/Elricriven/co2data> (Wang et al., 2022).

422

423 **6 Conclusions**

424 Based on the machine learning method, we reconstructed the sea surface $p\text{CO}_2$ fields in the SCS with high spatial resolution
425 ($0.05^\circ \times 0.05^\circ$) over the last two decades (2003-2020) by calculating the statistical relationship between the underway observational
426 $p\text{CO}_2$ data and remote sensing data. The machine learning method used in this study was facilitated by the EOF method, because



427 the latter can provide spatial constraints for the data reconstruction. In addition to the typical machine learning performance
428 metrics, we present a novel uncertainty calculation method that incorporates the bias of both the reconstruction and the sensitivity
429 of reconstructed models to its features. This method effectively shows the spatiotemporal patterns of bias, and makes up for the
430 spatial representation of the typical performance metrics.

431 We validate our reconstruction with three independent testing datasets, and the results show that the bias between our
432 reconstruction and observational $p\text{CO}_2$ data in the SCS is relatively small (about $10 \mu\text{atm}$). Our reconstruction successfully shows
433 the main features of the spatial and temporal patterns of $p\text{CO}_2$ in the SCS, indicating that we can use these reconstructed data to
434 further analyse the effect of meso-microscale processes (e.g., the Pearl River plume, and CCC) on sea surface $p\text{CO}_2$ in the SCS.

435 We divided the SCS into five sub-regions and separately calculated the deseasonalized long term trend of $p\text{CO}_2$ in each subregion,
436 and compared them with the long-term trend of atmospheric $p\text{CO}_2$. Our results show that the reconstructed data are consistent
437 with those of observational data. Moreover, the strength of the CO_2 sink in the northern SCS shows an increasing trend, whereas
438 $p\text{CO}_2$ trends in other subregions are essentially the same as that of atmospheric $p\text{CO}_2$.

439 This high spatiotemporal resolution of sea surface $p\text{CO}_2$ data is helpful to clarify the controlling factors of $p\text{CO}_2$ change in the
440 SCS and may be useful to predict changes of CO_2 source or sink patterns in this system.

441

442 **Author contribution**

443 Minhan Dai conceptualized and directed the field program of in situ observations. Xianghui Guo and Yi Xu participated in the in
444 situ data collection. Yan Bai provided the remote sensing-derived $p\text{CO}_2$ data. Minhan Dai, Guizhi Wang and Zhixuan Wang
445 developed the reconstruction method, wrote the codes, analyzed the data, and plotted the figures. Zhixuan Wang wrote the
446 manuscript. Minhan Dai and Xianghui Guo contributed to the writing, editing and revision of the original manuscript.

447

448 **Competing interests**

449 The authors declare that they have no conflict of interest.

450

451 **Acknowledgements**

452 We thank the support of the National Natural Science Foundation of China (grant No. 42188102, 42141001, and 41890800), and
453 the National Basic Research Program of China (973 Program, grant No.2015CB954000).

454

455 **References**

456 Bai, Y., Cai, W., He, X., Zhai, W., Pan, D., Dai, M., and Yu, P.: A mechanistic semi-analytical method for remotely sensing sea
457 surface $p\text{CO}_2$ in river-dominated coastal oceans: A case study from the East China Sea, *J. Geophys. Res.: Oceans*, 120,



- 458 2331-2349, 2015.
- 459 Bakker, D., Pfeil, B., Landa, C., Metzl, N., and Xu, S.: A multi-decade record of high-quality $f\text{CO}_2$ data in version 3 of the Surface
460 Ocean CO₂ Atlas (SOCAT), *Earth Syst. Sci. Data*, 8, 383-413, 2016.
- 461 Borges, A. V., Delille, B., and Frankignoulle, M.: Budgeting sinks and sources of CO₂ in the coastal ocean: Diversity of
462 ecosystems counts, *Geophys. Res. Lett.*, 32, L14601, 2005.
- 463 Cao, Z. and Dai, M.: Shallow-depth CaCO₃ dissolution: Evidence from excess calcium in the South China Sea and its export to
464 the Pacific Ocean, *Global Biogeochem. Cy.*, 25, GB2019, 2011.
- 465 Cao, Z., Dai, M., Zheng, N., Wang, D., Li, Q., Zhai, W., Meng, F., and Gan, J.: Dynamics of the carbonate system in a large
466 continental shelf system under the influence of both a river plume and coastal upwelling, *J. Geophys. Res.: Biogeo.*, 116,
467 G02010, 2011.
- 468 Cao, Z., Yang, W., Zhao, Y., Guo, X., Yin, Z., Du, C., Zhao, H., and Dai, M.: Diagnosis of CO₂ dynamics and fluxes in global
469 coastal oceans, *Natl. Sci. Rev.*, 7, 786-797, 2020.
- 470 Chen, C. and Borges, A. V.: Reconciling opposing views on carbon cycling in the coastal ocean: Continental shelves as sinks and
471 near-shore ecosystems as sources of atmospheric CO₂, *Deep-Sea Res. I*, 56, 578-590, 2009.
- 472 Chen, C., Lai, Z., Beardsley, R. C., Xu, Q., Lin, H., and Viet, N. T.: Current separation and upwelling over the southeast shelf
473 of Vietnam in the South China Sea, *J. Geophys. Res.: Oceans*, 117, C03033, 2012.
- 474 Chen, F., Cai, W. J., Benitez-Nelson, C., and Wang, Y.: Sea surface $p\text{CO}_2$ -SST relationships across a cold-core cyclonic eddy:
475 Implications for understanding regional variability and air-sea gas exchange, *Geophys. Res. Lett.*, 34, 265-278, 2007.
- 476 Cheng, C., Xu, P. F., Cheng, H., Ding, Y., Zheng, J., Ge, T., and Xu, J.: Ensemble learning approach based on stacking for
477 unmanned surface vehicle's dynamics. *Ocean Eng.*, 207, 107388, 2020.
- 478 Dai, M. H., Cao, Z., Guo, X., Zhai, W., Liu, Z., Yin, Q., Xu, Y., Gan, J., Hu, J., and Du, C.: Why are some marginal seas sources
479 of atmospheric CO₂?, *Geophys. Res. Lett.*, 40, 2154-2158, 2013.
- 480 Dai, M., Gan, J., Han, A., Kung, H., and Yin, Z.: "Physical Dynamics and Biogeochemistry of the Pearl River Plume" in
481 *Biogeochemical Dynamics at Large River-Coastal Interfaces*. Eds. T. Bianchi, M. Allison and W. J. Cai (Cambridge University
482 Press, Cambridge), 321-352, 2014.
- 483 Dai, M., J. Su, Zhao, Y., Hofmann, E. E., Cao, Z., Cai, W., Gan, J., Lacroix, F., Laruelle, G., Meng, F., Müller, J., Regnier, P., Wang,
484 G., and Wang, Z.: Carbon fluxes in the coastal ocean: Synthesis, boundary processes and future trends, *Annu. Rev. Earth Pl. Sc.*,
485 50, 593-626, 2022.
- 486 Du, C., Liu, Z., Dai, M., Kao, S. J., and Li, Y.: Impact of the Kuroshio intrusion on the nutrient inventory in the upper northern
487 South China Sea: insights from an isopycnal mixing model, *Biogeosciences*, 10, 6419-6432, 2013.
- 488 Dong, L., Su, J. Wong, L. Cao, Z. and Chen, J.: Seasonal variation and dynamics of the Pearl River plume, *Cont. Shelf*



- 489 Res., 24, 1761-1777, 2004.
- 490 Fassbender, A. J., Rodgers, K. B., Palevsky, H. I., and Sabine, C. L.: Seasonal Asymmetry in the Evolution of Surface Ocean
491 $p\text{CO}_2$ and pH Thermodynamic Drivers and the Influence on Sea-Air CO_2 Flux, *Global Biogeochem. Cy.*, 32, 1476-1497, 2018.
- 492 Fay, A., Gregor, L., Landschützer, P., McKinley, G., Gruber, N., Gehlen, M., Iida, Y., Laruelle, G., Rödenbeck, C., Roobaert, A.,
493 and Zeng, J.: SeaFlux: harmonization of air-sea CO_2 fluxes from surface $p\text{CO}_2$ data products using a standardized approach,
494 *Earth Syst. Sci. Data*, 13, 4693-4710, 2021
- 495 Friedlingstein, P., Jones, M. W., O'Sullivan, M., Andrew, R. M., Hauck, J., Peters, G. P., Peters, W., Pongratz, J., Sitch, S., Le
496 Quéré, C., Bakker, D. C. E., Canadell, J. G., Ciais, P., Jackson, R. B., Anthoni, P., Barbero, L., Bastos, A., Bastrikov, V., Becker,
497 M., Bopp, L., Buitenhuis, E., Chandra, N., Chevallier, F., Chini, L. P., Currie, K. I., Feely, R. A., Gehlen, M., Gilfillan, D.,
498 Gkritzalis, T., Goll, D. S., Gruber, N., Gutekunst, S., Harris, I., Haverd, V., Houghton, R. A., Hurtt, G., Ilyina, T., Jain, A. K.,
499 Joetzier, E., Kaplan, J. O., Kato, E., Klein Goldewijk, K., Korsbakken, J. I., Landschützer, P., Lauvset, S. K., Lefèvre, N.,
500 Lenton, A., Lienert, S., Lombardozi, D., Marland, G., McGuire, P. C., Melton, J. R., Metzl, N., Munro, D. R., Nabel, J. E. M.
501 S., Nakaoka, S.-I., Neill, C., Omar, A. M., Ono, T., Peregon, A., Pierrot, D., Poulter, B., Rehder, G., Resplandy, L., Robertson,
502 E., Rödenbeck, C., Séférian, R., Schwinger, J., Smith, N., Tans, P. P., Tian, H., Tilbrook, B., Tubiello, F. N., van der Werf, G. R.,
503 Wiltshire, A. J., and Zaehle, S.: Global Carbon Budget 2019, *Earth Syst. Sci. Data*, 11, 1783-1838, 2019.
- 504 Gan, J., Li, H., Curchitser, E. N., and Haidvogel, D. B.: Modeling South China sea circulation: Response to seasonal forcing
505 regimes, *J. Geophys. Res.: Oceans*, 111, C06034, 2006.
- 506 Gan, J., Li, L., Wang, D., and Guo, X.: Interaction of a river plume with coastal upwelling in the northeastern South China Sea,
507 *Cont. Shelf Res.*, 29, 728-740, 2009.
- 508 Gan, J., Lu, Z., Dai, M., Cheung, A. Y. Y., Liu, H., and Harrison, P.: Biological response to intensified upwelling and to a river
509 plume in the northeastern South China Sea: A modeling study, *J. Geophys. Res.: Oceans*, 115, C09001, 2010.
- 510 Guo, X. and Wong, G.: Carbonate chemistry in the Northern South China Sea shelf-sea in June 2010, *Deep Sea Res. II*, 117,
511 119-130, 2015.
- 512 Han, A. Q., Dai, M. H., Gan, J. P., Kao, S. J., Zhao, X. Z., Jan, S., Li, Q., Lin, H., Chen, C. T. A., and Wang, L.: Inter-shelf nutrient
513 transport from the East China Sea as a major nutrient source supporting winter primary production on the northeast South China
514 Sea shelf, *Biogeosciences*, 10, 8159-8170, 2013.
- 515 Hu, J., Kawamura, H., Li, C., Hong, H., and Jiang, Y.: Review on current and seawater volume transport through the Taiwan Strait,
516 *J. Oceanogr.*, 66, 591-610, 2010.
- 517 Jones, S. D., Quéré, C., and Rödenbeck, C.: Spatial decorrelation lengths of surface ocean $f\text{CO}_2$ results in NetCDF format, *Global*
518 *Biogeochem. Cy.*, 26, GB2042, 2014.
- 519 Jo, Y., Dai, M., Zhai, W., Yan, X., and Shang, S.: On the Variations of Sea Surface $p\text{CO}_2$ in the Northern South China Sea - A



- 520 Remote Sensing Based Neural Network Approach, *J. Geophys. Res.: Oceans*, 117, C08022, 2012.
- 521 Landschützer, P., Gruber, N., and Bakker, D.: Decadal variations and trends of the global ocean carbon sink, *Global Biogeochem.*
522 *Cy.*, 30, 1396-1417, 2016.
- 523 Landschützer, P., Gruber, N., and Bakker, D. C. E.: An updated observation-based global monthly gridded sea surface $p\text{CO}_2$ and
524 air-sea CO_2 flux product from 1982 through 2015 and its monthly climatology, Dataset, 2017.
- 525 Laruelle, G., Lauerwald, R., Pfeil, B., and Regnier, P.: Regionalized global budget of the CO_2 exchange at the air-water interface
526 in continental shelf seas, *Global Biogeochem. Cy.*, 28, 1199-1214, 2015.
- 527 Landschützer, P., Bakker, D. C. E., Gruber, N., and Schuster, U.: Recent variability of the global ocean carbon sink, *Global*
528 *Biogeochem. Cy.*, 28, 927-949, 2014.
- 529 Lefèvre, N., Watson, A., and Waston, A.: A comparison of multiple regression and neural network techniques for mapping in situ
530 $p\text{CO}_2$ data, *Tellus B*, 57, 375-384, 2005.
- 531 Lefèvre, N., Watson, A. J., and Watson, A. R.: A comparison of multiple regression and neural network techniques for mapping in
532 situ $p\text{CO}_2$ data, *Tellus B: Chemical and Physical Meteorology, Dataset*, 2017.
- 533 Li, Y., Xie, P., Tang, Z., Jiang, T., and Qi, P.: SVM-Based Sea-Surface Small Target Detection: A False-Alarm-Rate-Controllable
534 Approach, *IEEE Geosc.i Remote S.*, 16, 1225-1229, 2019.
- 535 Li, H., Wiesner, M. G., Chen, J., Lin, Z., Zhang, J., and Ran, L.: Long-term variation of mesopelagic biogenic flux in the central
536 South China Sea: Impact of monsoonal seasonality and mesoscale eddy, *Deep Sea Res. I*, 126, 62-72, 2017.
- 537 Li, Q., Guo, X., Zhai, W., Xu, Y., Dai, M.: Partial pressure of CO_2 and air-sea CO_2 fluxes in the South China Sea: Synthesis of an
538 18-year dataset, *Prog. Oceanogr.*, 182, 102272, 2020.
- 539 Luo, X., Hao, W., Zhe, L., and Liang, Z.: Seasonal variability of air-sea CO_2 fluxes in the Yellow and East China Seas: A case
540 study of continental shelf sea carbon cycle model, *Cont. Shelf Res.*, 107, 69-78, 2015.
- 541 Mongwe, N. P., Chang, N., and Monteiro, P.: The seasonal cycle as a mode to diagnose biases in modelled CO_2 fluxes in the
542 Southern Ocean, *Ocean Model.*, 106, 90-103, 2016.
- 543 Park, J. H.: Effects of Kuroshio intrusions on nonlinear internal waves in the South China Sea during winter, *J. Geophys. Res.:*
544 *Oceans*, 118, 7081-7094, 2013.
- 545 Qin, H., Chen, G., Wang, W., Wang, D., and Zeng, L.: Validation and application of MODIS-derived SST in the South China Sea,
546 *Int. J. Remote Sens.*, 35, 4315-4328, 2014.
- 547 Rödenbeck, C., Bakker, D. C. E., Gruber, N., Iida, Y., Jacobson, A. R., Jones, S., Landschützer, P., Metzl, N., Nakaoka, S., and
548 Olsen, A.: Data-based estimates of the ocean carbon sink variability-first results of the Surface Ocean $p\text{CO}_2$ Mapping
549 intercomparison (SOCOM), *Biogeosciences*, 12, 14-49, 2015.
- 550 Sheu, D. D., Chou, W. C., Wei, C. L., Hou, W. P., Wong, G., and Hsu, C. W.: Influence of El Niño the sea-to-air CO_2 flux at the



- 551 SEATs time-series site, northern South China Sea, *J. Geophys. Res.: Oceans*, 115, C10021, 2010.
- 552 Tahata, M., Sawaki, Y., Ueno, Y., Nishizawa, M., Yoshida, N., Ebisuzaki, T., Komiya, T., and Maruyama, S.: Three-step
553 modernization of the ocean: Modeling of carbon cycles and the revolution of ecological systems in the Ediacaran/Cambrian
554 periods, *Geosci. Front.*, 6, 121-136, 2015.
- 555 Telszewski, M., Chazottes, A., Schuster, U., Watson, A. J., and Wanninkhof, R.: Estimating the monthly $p\text{CO}_2$ distribution in the
556 North Atlantic using a self-organizing neural network, *Biogeosciences*, 6, 1405-1421, 2009.
- 557 Wang, G., Shen, S. S. P., Chen, Y., Bai, Y., Qin, H., Wang, Z., Chen, B., Guo, X., and Dai, M.: Feasibility of reconstructing the
558 basin-scale sea surface partial pressure of carbon dioxide from sparse in situ observations over the South China Sea, *Earth Syst.*
559 *Sci. Data*, 13, 1403-1417, 2021.
- 560 Wanninkhof, R., Park, G. H., Takahashi, T., Sweeney, C., Feely, R., Nojiri, Y., Gruber, N., Doney, S. C., Mckinley, G. A., and
561 Lenton, A.: Global ocean carbon uptake: magnitude, variability and trend, *Biogeosciences*, 10, 1983-2000, 2013
- 562 Wang, Z., Wang, G., Guo, X., Bai, Y., Xu, Y., and Dai, M.: Spatial reconstruction of long-term (2003-2020) sea surface $p\text{CO}_2$
563 in the South China Sea using a machine learning based regression method aided by empirical orthogonal function
564 analysis, Github, <https://github.com/Elricriven/co2data>.
- 565 Xu, X., Zang, K., Zhao, H., Zheng, N., Huo, C., and Wang, J.: Monthly CO_2 at A4HDYD station in a productive shallow marginal
566 sea (Yellow Sea) with a seasonal thermocline: Controlling processes, *J. Marine Syst.*, 159, 89-99, 2016.
- 567 Jo, Y., Dai, M., Zhai, W., Yan, X., and Shang, S.: On the variations of sea surface $p\text{CO}_2$ in the northern South China Sea: A remote
568 sensing based neural network approach, *J. Geophys. Res.: Oceans*, 117, C08022, 2012.
- 569 Yang, W., Guo, X., Cao, Z., Wang, L., Guo, L., Huang, T., Li, Y., Xu, Y., Gan, J., and Dai, M.: Seasonal dynamics of the carbonate
570 system under complex circulation schemes on a large continental shelf: The northern South China Sea, *Prog Oceanogr.*, 197,
571 1026-1045, 2021.
- 572 Yu, Z., Shang, S., Zhai, W., and Dai, M.: Satellite-derived surface water $p\text{CO}_2$ and air-sea CO_2 fluxes in the northern South China
573 Sea in summer, *Prog. Nat. Sci.*, 19, 775-779, 2009.
- 574 Zeng, J., Matsunaga, T., Saigusa, N., Shirai, T., Nakaoka, S. I., and Tan, Z. H.: Technical note: Evaluation of three machine
575 learning models for surface ocean CO_2 mapping, *Ocean Sci.*, 13, 303-313, 2017.
- 576 Zhai, W., Dai, M., Cai, W. J., Wang, Y., and Hong, H.: The partial pressure of carbon dioxide and air-sea fluxes in the northern
577 South China Sea in spring, summer and fall, *Mar. Chem.*, 96, 87-97, 2005.
- 578 Zhai, W. D., Dai, M. H., Chen, B. S., Guo, X. H., Li, Q., Shang, S. L., Zhang, C. Y., Cai, W. J., and Wang, D. X.: Seasonal
579 variations of sea-air CO_2 fluxes in the largest tropical marginal sea (South China Sea) based on multiple-year underway
580 measurements, *Biogeosciences*, 10, 7775-7791, 2013.
- 581 Zhang, C., Hu, C., Shang, S., Müller-Karger, F., Yan, L., Dai, M., Huang, B., Ning, X., and Hong, H.: Bridging between SeaWiFS



- 582 and MODIS for continuity of chlorophyll-a concentration assessments off Southeastern China, *Remote Sens. Environ.*, 102,
583 250-263, 2006.
- 584 Zhan, Y., Zhang, H., Li, J., and Li, G.: Prediction Method for Ocean Wave Height Based on Stacking Ensemble Learning Model. *J.*
585 *Mar. Sci. Eng.*, 10, 1150, 2022.
- 586 Zhu, Y., Shang, S., Zhai, W., and Dai, M.: Satellite-derived surface water $p\text{CO}_2$ and air-sea CO_2 fluxes in the northern South China
587 Sea in summer, *Prog. Nat. Sci.*, 19, 775-779, 2009.