

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

Zhixuan Wang<sup>1</sup>, Guizhi Wang<sup>1,2</sup>, Xianghui Guo<sup>1</sup>, Yan Bai<sup>3</sup>, Yi Xu<sup>1</sup> and Minhan Dai<sup>1,\*</sup>

<sup>1</sup>State Key Laboratory of Marine Environmental Science and College of Ocean and Earth Sciences, Xiamen University, Xiamen, 361102, China

<sup>2</sup>Fujian Provincial Key Laboratory for Coastal Ecology and Environmental Studies, Xiamen University, Xiamen, 361102, China

<sup>3</sup>State Key Laboratory of Satellite Ocean Environment Dynamics, Second Institute of Oceanography, State Oceanic Administration, Hangzhou, 310012, China

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

**Abstract.** The South China Sea (SCS) is the largest marginal sea in the North Pacific Ocean, where intensive field observations including mappings of the sea-surface partial pressure of  $\text{CO}_2$  ( $p\text{CO}_2$ ) have been conducted over the last two decades. It is one of the most studied marginal seas in terms of carbon cycling, and could thus be a model system for marginal sea carbon research. However, the cruise-based sea surface  $p\text{CO}_2$  datasets are still temporally and spatially sparse. Using a machine learning-based method facilitated by empirical orthogonal function (EOF) analysis, this study provides a reconstructed dataset of the monthly sea surface  $p\text{CO}_2$  in the SCS with a reasonably high spatial resolution ( $0.05^\circ \times 0.05^\circ$ ) and temporal coverage between 2003 and 2020.

The data input to our reconstructed model includes remote sensing-derived sea surface salinity, sea surface temperature, and chlorophyll, the spatial pattern of  $p\text{CO}_2$  constrained by EOF, atmospheric  $p\text{CO}_2$ , and time-labels (month). We validated our reconstruction with three independent testing datasets that are not involved in the model training. Among them, Test 1 includes 10% of our *in situ* data, Test 2 contains four independent *in situ* datasets corresponding to the four seasons, and Test 3 is an *in situ* monthly dataset available from 2003–2019 at the South East Asia Time-Series (SEATs) station located in the northern basin of the SCS. Our Test 1 validation demonstrated that the reconstructed  $p\text{CO}_2$  field successfully simulated the spatial and temporal patterns of sea surface  $p\text{CO}_2$  observations. The root-mean-square error (RMSE) between our reconstructed data and *in situ* data in Test 1 averaged  $\sim 10 \mu\text{atm}$ , which is much smaller (by  $\sim 50\%$ ) than that between the remote sensing-derived data and *in situ* data. Test 2 verified the accuracy of our retrieval algorithm in months lacking observations, showing a relatively small bias (RMSE:  $\sim 8 \mu\text{atm}$ ). Test 3 evaluated the accuracy of the reconstructed long-term trend, showing that at the SEATs Station, the difference between the reconstructed  $p\text{CO}_2$  and *in situ* data ranged from  $-10$  to  $4 \mu\text{atm}$  ( $-2.5\%$  to  $1\%$ ). In addition to the typical machine learning performance metrics, we assessed the uncertainty resulting from reconstruction bias and its feature sensitivity. These validations

删除[Author]: ocean

删除[Author]: of

删除[Author]: and constraints of spatiotemporal modes of

删除[Author]: The SCS

删除[Author]: and

删除[Author]: .

删除[Author]: mapping sea surface  $p\text{CO}_2$  of this region is

删除[Author]: which

删除[Author]: The South China Sea (SCS) is the largest

删除[Author]: datasets of

删除[Author]: incomplete

删除[Author]: capable of constraining the spatiality

删除[Author]: and selecting the remote sensing derived

删除[Author]: input

删除[Author]: in

删除[Author]: O

删除[Author]: ion was initiated by using

删除[Author]: s

删除[Author]: the

删除[Author]:

删除[Author]: data (, which include

删除[Author]: and

删除[Author]: s

删除[Author]: calculated

删除[Author]: )

删除[Author]:

删除[Author]: as inputs data

删除[Author]: ,

删除[Author]: (it indicates that this part of the data which

删除[Author]: is completely

删除[Author]: un

删除[Author]: )

删除[Author]: where,

删除[Author]: EST

and uncertainty analyses strongly suggest that our reconstruction effectively captures the main spatial and temporal features of sea surface  $p\text{CO}_2$  distributions in the 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 differing physico-biogeochemical characteristics. We show that mesoscale processes such as the Pearl River plume and China Coastal Currents significantly impact sea surface  $p\text{CO}_2$  in the SCS during different seasons. While the SCS is overall a weak source of atmospheric  $\text{CO}_2$ , the northern SCS acts as a sink, showing a trend of increasing strength over the past two decades.

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

## 1 Introduction

The ocean possesses a large portion of the global capacity for atmospheric carbon dioxide ( $\text{CO}_2$ ) sequestration, annually mitigating 22%–26% of the anthropogenic  $\text{CO}_2$  emissions associated with fossil fuel burning and land use changes over the period from 2012–2021 (Friedlingstein et al., 2022). Ocean margins are an essential part of the land-ocean continuum, representing a particularly challenging regime to study (e.g., Chen and Borges, 2009; Dai et al. 2022; Laruelle et al., 2014), as they are often characterized by large spatial and temporal variations in air-sea  $\text{CO}_2$  fluxes that lead to larger uncertainties in their overall estimation and predictions than those made in the open ocean (Dai et al., 2013, 2022; Cao et al., 2020; Laruelle et al., 2014; Chen and Borges, 2009 and the references therein). Limited spatiotemporal coverage of in situ observations is a large source of these uncertainties.

In recent years, many studies have used numerical models or data-based approaches to improve estimates of the partial pressure of carbon dioxide ( $p\text{CO}_2$ ) at the sea surface and the accuracy of the global carbon budget for periods and regions with poor coverage of in situ data (e.g., Rödenbeck et al., 2015; Wanninkhof et al., 2013). Numerical models can successfully quantify the generally increasing trend in oceanic  $p\text{CO}_2$  and simulate some critical carbon cycling processes (e.g., net ecosystem production), but still suffer from regional and seasonal differences in their estimates of ocean carbonate parameters (e.g., Luo et al., 2015; Mongwe et al., 2016; Tahata et al., 2015; Wanninkhof et al., 2013). Thus, data-based approaches, which typically apply statistical interpolation and regression methods, have become an important complement to numerical models (e.g., Jones et al., 2014; Lefèvre et al., 2005; Landschützer et al., 2014, 2017; Telszewski et al., 2009). Statistical interpolation improves the spatial coverage of in situ data, but does not work for periods where in situ data are unavailable. Regression methods allow mapping of the relationships between in situ  $p\text{CO}_2$  data and other parameters that may drive changes in surface ocean  $p\text{CO}_2$ , and then the extrapolation of this relationship to improve estimates of the spatiotemporal distribution of  $p\text{CO}_2$ . Machine learning methods and remote sensing-derived products (as proxy variables in regression methods) have aided the development of data-based methods (Rödenbeck et al., 2015; Bakker et al., 2016), and can improve the model results for the oceanic carbonate system by numerical

删除[Author]: analysis

删除[Author]: in both the spatial and temporal patterns

删除[Author]: much

删除[Author]: and

删除[Author]: essd

删除[Author]: the

删除[Author]: s

删除[Author]: during

删除[Author]: 1960

删除[Author]: 19

删除[Author]: 0

删除[Author]: However, it remains largely unknown whether

删除[Author]: ,

删除[admin]: and despite occupying only 7% of the surface

删除[Author]: and

删除[admin]:

删除[Author]: . This large uncertainty is primarily attributed

删除[Author]: bilities

删除[Author]: y

删除[Author]: of

删除[Author]: even

删除[Author]: y

删除[Author]: prediction

删除[Author]: those

删除[Author]: occurring

删除[Author]: al data

删除[Author]: n important

删除[Author]: sea surface

删除[Author]: partial pressure

删除[Author]:  $\text{CO}_2$  distribution

删除[Author]: observational

删除[admin]: ocean

删除[Author]: of performance

删除[Author]:  $\text{CO}_2$

assimilation methods. Consequently, machine learning has increasingly become a routine approach for reconstructing sea surface  $p\text{CO}_2$  in open ocean regimes (e.g., Zeng et al., 2017; Li et al., 2019); however, it remains challenging to extend this method to ocean margins, which are more dynamic in both time and space

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 extensive field observations of sea surface  $p\text{CO}_2$  have been conducted in the SCS over the past two decades, their spatial and temporal coverage is still limited with respect to coverage of different physical-biogeochemical domains and sub-seasonal time scales (e.g., Guo et al., 2015; Li et al., 2020; Zhai et al., 2005; Zhai et al., 2013). Therefore, there is a strong need for improved surface water  $p\text{CO}_2$  coverage in the SCS to constrain air-sea  $\text{CO}_2$  fluxes and improve initial conditions of numerical models.

Moreover, reasonably high 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 long-term changes.

Zhu et al. (2009) presented an empirical approach to estimate sea surface  $p\text{CO}_2$  in the northern SCS using remote sensing-derived (RS-derived) data, including sea surface temperature (SST) and chlorophyll  $a$  (Chl  $a$ ). Their

reconstructed  $p\text{CO}_2$  data were generally consistent with the in situ data. However, uncertainties remained large, primarily caused by limited in situ data from only two summer cruises in their study. Jo et al. (2012) developed a neural network-based algorithm

using SST and Chl  $a$  to estimate sea surface  $p\text{CO}_2$  in the northern SCS. In their study, in situ sea surface  $p\text{CO}_2$  data were collected from three cruises during May 2001 and February and July 2004. The reconstruction also suffered a relatively large bias, (Wang et

al., 2021). Bai et al. (2015) employed a ‘mechanic semi-analytical algorithm (MeSAA)’ to estimate satellite remote sensing-derived sea surface  $p\text{CO}_2$  in the East China Sea from 2000–2014, and then expanded the application of this algorithm to

estimate sea surface  $p\text{CO}_2$  for the whole China Seas region including the South China Sea. These authors explained that their MeSAA did not fully account for some localized processes, which resulted in a RMSE of about  $45 \mu\text{atm}$  for the SCS (Wang et al.,

2021). Yu et al. (2022) subsequently used a non-linear regression method to develop a retrieval algorithm for seawater  $p\text{CO}_2$  in the

China Seas, and the RS-derived  $p\text{CO}_2$  data from 2003-2018 were provided by the Sat $\text{CO}_2$  platform (www.Sat $\text{CO}_2$ .com). In this

retrieval algorithm, the input parameters included sea surface temperature, Chl  $a$  concentrations, remote sensing reflectance at

three bands (Rrs412, 443, 488 nm), the temperature anomaly in the longitudinal direction, and the theoretical thermodynamic

background  $p\text{CO}_2$  under the corresponding SST. Although the RMSE associated with the RS-derived  $p\text{CO}_2$  product was relatively

large ( $21.1 \mu\text{atm}$ ), it successfully showed the major spatial patterns of sea surface  $p\text{CO}_2$  in the China Seas (Yu et al., 2022).

To take advantage of both the high spatiotemporal resolution of the RS-derived  $p\text{CO}_2$  data, and the accuracy of the in situ data,

Wang et al. (2021) reconstructed a basin-scale sea surface  $p\text{CO}_2$  dataset in the SCS during summer using an empirical orthogonal

function (EOF) based on a multi-linear regression method. They demonstrated that the spatial modes of RS-derived data

calculated using the EOF can effectively provide spatial constraints on the data reconstruction, and thus this approach is adopted

in this study. However, the reconstructed results may still be subject to bias when the standard deviation of spatial in situ data is

删除[Author]: Thus

删除[Author]: been

删除[Author]: widely used forin

删除[Author]: the

删除[Author]: on

删除[Author]: of

删除[Author]: for the global ocean

删除[Author]: (refs?)

删除[Author]: .

删除[Author]: , hH

删除[Author]: still

删除[admin]: marginal seas

删除[Author]: featuringe

删除[Author]: changes in both time and spacespatially an ...

删除[Author]: have been conducted

删除[Author]: in

删除[Rick Smith]: in different

删除[Rick Smith]: of the SCS

删除[Rick Smith]: at

删除[Author]: clear

删除[Rick Smith]: to achieve

删除[Rick Smith]: with a highest spatiotemporal resolutic ...

删除[Rick Smith]: in the SCS

删除[Author]: and

删除[Rick Smith]: so as tohelp

删除[Author]: develop

删除[Author]: d

删除[Author]: a

删除[Author]: the

删除[Author]: that

删除[Author]: d

删除[Author]: in summer

删除[Author]: satellite

删除[Author]: -

92 relatively large because of the influence of outliers (Wang et al., 2021). Therefore, many studies have used machine  
93 learning-based regression methods to reduce the influence of outliers in open ocean areas, and have achieved a RMSE of  
94 <math>17 \mu\text{atm}</math> in most cases (e.g., Zeng et al., 2017; Li et al., 2019).

95 Building on the ability of the EOF method to significantly improve reconstructions in terms of spatial patterns and accuracy  
96 (Wang et al., 2021), we developed a machine learning-based regression method facilitated by the EOF to fully resolve the  
97 long-term spatial distribution of sea surface  $p\text{CO}_2$  at a resolution of  $0.05^\circ \times 0.05^\circ$  in the SCS. Our reconstructed model uses input  
98 data that includes remote sensing-derived sea surface salinity, sea surface temperature, and  $\text{Chl } a$ , the spatial pattern of  $p\text{CO}_2$   
99 constrained by the EOF, atmospheric  $p\text{CO}_2$ , and time labels (month). In addition to assessing typical machine learning  
100 performance metrics, we evaluated the uncertainty resulting from the bias of the reconstruction and its sensitivity to the features.  
101

## 102 2 Study site and data sources

### 103 2.1 Study area

104 The SCS, located in the northwestern Pacific, is a semi-enclosed marginal sea with a maximum water depth of ca. 4700 m (e.g.,  
105 Gan et al., 2006, 2010). The rhombus-shaped deep-water basin, with a southwest-northeast direction, accounts for about half of  
106 the total area of the SCS (Figure 1). Largely modulated by the Asian monsoon and topography, the SCS exhibits seasonally  
107 varying surface circulation, river inputs, and upwelling. The circulation of the upper layer shows a large cyclonic circulation  
108 structure in winter (Figure 1), while in summer it exhibits an anticyclonic circulation structure (Figure 1; Hu et al. 2010). In the  
109 northern SCS, the Pearl River discharges into the SCS with an annual freshwater input of  $3.26 \times 10^{11} \text{ m}^3$  (e.g., Dong et al., 2004;  
110 Dai et al., 2014). The area influenced by the Pearl River plume may extend southeastward to a few hundred kilometers from the  
111 estuary in summer because of the monsoonal wind stress (Dai et al., 2014). The northern and western coastal regions of the SCS  
112 feature summer coastal upwelling, such as the Eastern Guangdong and Qiongdong upwelling systems in the northern SCS and the  
113 Vietnam upwelling 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  
114 al., 2020). These seasonal changes of sea surface circulation lead to strong seasonal characteristics of sea surface  $p\text{CO}_2$  in the  
115 SCS.  
116

设置格式[admin]: 非突出显示

删除[Author]: the

删除[Author]: for

删除[Rick Smith]: with

删除[Author]: of

删除[Rick Smith]: up

删除[Rick Smith]: that

删除[Rick Smith]: d

删除[Rick Smith]: the

删除[Author]: ,

删除[Rick Smith]: Aand the input data in

删除[Rick Smith]: o

删除[Rick Smith]: include

删除[admin]: chlorophyll

删除[Author]: the

删除[Rick Smith]: assessed

删除[Author]: we present a novel uncertainty calculation

删除[XHGuo]: sea basin

删除[admin]: has

删除[Author]: .

删除[admin]: The oceanography of the SCS is l

删除[Rick Smith]: the

删除[admin]: thus

删除[admin]: ing

删除[admin]: Forced by the northeast winds in winter, t

删除[Author]: red solid line in

删除[Author]: .

删除[admin]: forced by southwest winds

删除[Author]: red dashed line in

删除[Author]: .

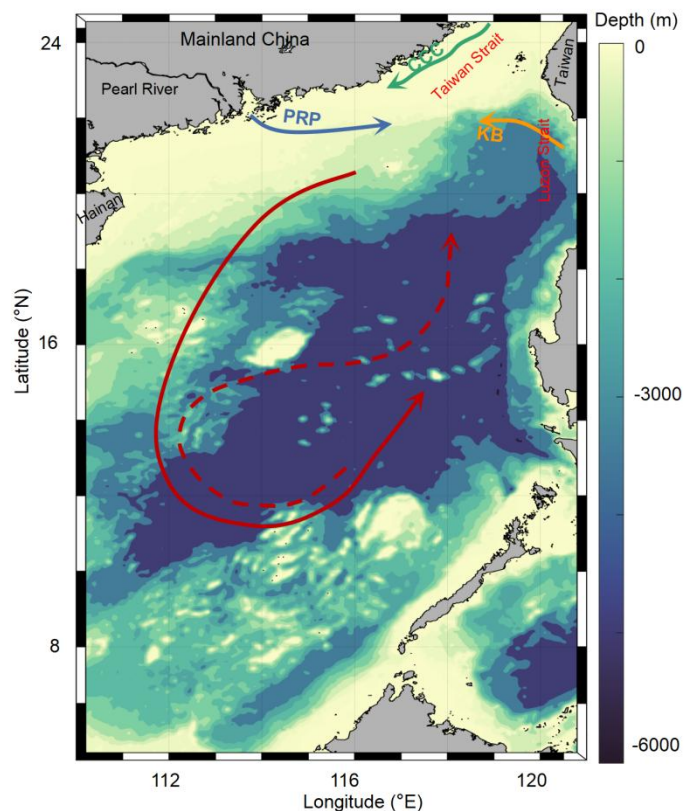
删除[Author]: P

删除[Author]: river

删除[admin]: also

删除[Rick Smith]: in summer





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

The SCS is subject to dynamic water exchanges with the East China Sea via the Taiwan Strait and the Western Pacific via the Luzon Strait (Fig. 1). In winter, driven by the winter monsoon, the China Coastal Current (CCC, green line in Fig. 1; Han et al., 2013; Yang et al., 2022) flows south along the Chinese mainland through the Taiwan Strait, and occupies the northern SCS with cold, fresh, nutrient-rich waters. The strong northeast winds in winter also slow down the western boundary ocean current, forcing the intrusion of Kuroshio water featuring high surface salinity and high total alkalinity, into the SCS via the Luzon Strait (orange line in Fig. 1; Du et al., 2013; Park, 2013; Yang et al., 2022). These water exchange processes increase the complexity of the spatial distribution of sea surface  $p\text{CO}_2$  in the SCS, which as a result, has strong seasonal characteristics and spatial variability.

## 2.2 Observational $p\text{CO}_2$ data

Data collected from field surveys during the study period 2003-2020 are summarized in Table 1. Most observations were made in July, with fewer observations made in March and December of each year. The rough sea state in the SCS in winter and early spring limited the field surveys during these seasons. Data collected from July 2000 to January 2018 were originally published by Li et al. (2020). The in situ  $p\text{CO}_2$  were collected from R/Vs *Dongfanghong-2*, *Tan Kah Kee (TKK)* (shown in Table 1). During the cruises, sea surface  $p\text{CO}_2$  was measured during the cruise. The measurements and data processing followed the SOCAT (Surface Ocean CO<sub>2</sub> Atlas) protocol (Li et al., 2020). More details of the data collection methods are provided in Li et al. (2020). The

删除[Rick Smith]: a  
 删除[Author]: gyre  
 删除[Rick Smith]: an  
 删除[Author]: gyre  
 设置格式[admin]: 图案: 清除(自动设置), 非突出显示  
 删除[admin]: is a semi-enclosed sea basin with  
 设置格式[admin]: 图案: 清除(自动设置), 非突出显示  
 删除[Author]: yellow  
 删除[admin]: , which shows  
 删除[Rick Smith]: . A  
 删除[admin]: s  
 删除[Rick Smith]: a resul  
 删除[Rick Smith]: ,  
 删除[admin]: A high-spatial-resolution sea surface  
 删除[Author]: carbon cycle  
 删除[admin]: .  
 删除[admin]: Most  
 删除[Rick Smith]: and  
 删除[Rick Smith]: were  
 删除[Rick Smith]:  
 删除[Zhixuan Wang]: , etc.  
 删除[admin]: .  
 删除[Author]: continuously  
 删除[Author]: methods  
 删除[Author]: those of  
 删除[Author]: , <http://www.socat.info/news.html>  
 删除[Author]: protocol  
 删除[Author]: T  
 删除[Author]: used in this study  
 删除[Rick Smith]: have been introduced

136 spatial coverage and frequency of the observations are shown in Figure 2, revealing pronounced seasonal changes across a large  
 137 spatial area. For example, the spatial coverage of the in situ data in spring and fall are relatively uniformly distributed, and the  
 138 south end of the spatial coverage reaches 5 °N in spring, whereas during other seasons the data are concentrated in the northern  
 139 and central regions of the SCS. In addition, only one observation was made in the basin area in winter, while the northern coastal  
 140 area was more frequently surveyed, especially in summer.

141 **Table 1. Summary of seasonal in situ data of sea surface pCO<sub>2</sub> in the South China Sea for the period 2003-2020 used in this**  
 142 **study.**

<u>Season</u>	<u>Spring</u>				<u>Summer</u>	
	<u>March</u>	<u>April</u>	<u>May</u>	<u>June</u>	<u>July</u>	<u>August</u>
<u>Cruise</u> <u>time</u>					<u>2004.07</u>	
		<u>2005.04</u>		<u>2006.06</u>	<u>2005.07</u>	
		<u>2008.04</u>	<u>2004.05</u>	<u>2016.06</u>	<u>2007.07</u>	<u>2007.08</u>
	<u>2004.03</u>	<u>2009.04</u>	<u>2011.05</u>	<u>2017.06*</u>	<u>2008.07</u>	<u>2008.08</u>
		<u>2012.04</u>	<u>2014.05</u>	<u>2019.06*</u>	<u>2009.07</u>	<u>2019.08*</u>
		<u>2020.04*</u>	<u>2020.05*</u>	<u>2020.06*</u>	<u>2012.07</u>	
					<u>2015.07*</u>	
					<u>2019.07*</u>	
<u>Season</u>	<u>Fall</u>			<u>Winter</u>		
	<u>September</u>	<u>October</u>	<u>November</u>	<u>December</u>	<u>January</u>	<u>February</u>
<u>Cruise</u> <u>time</u>	<u>2004.09</u>				<u>2009.01</u>	
	<u>2007.09</u>	<u>2003.10</u>	<u>2006.11</u>		<u>2010.01</u>	<u>2004.02</u>
	<u>2008.09</u>	<u>2006.10</u>	<u>2010.11</u>	<u>2006.12</u>	<u>2018.01</u>	<u>2006.02</u>
	<u>2020.09*</u>					
<u>Data</u> <u>source</u>			<u>Li et al. (2020)</u>			
			<u>*This study</u>			

删除[admin]: and  
 删除[Rick Smith]: show that there are  
 删除[Rick Smith]: and that the data cover  
 删除[Author]: observed  
 删除[Author]: that  
 删除[Author]: is  
 删除[Author]: the  
 删除[Author]: **observational data**

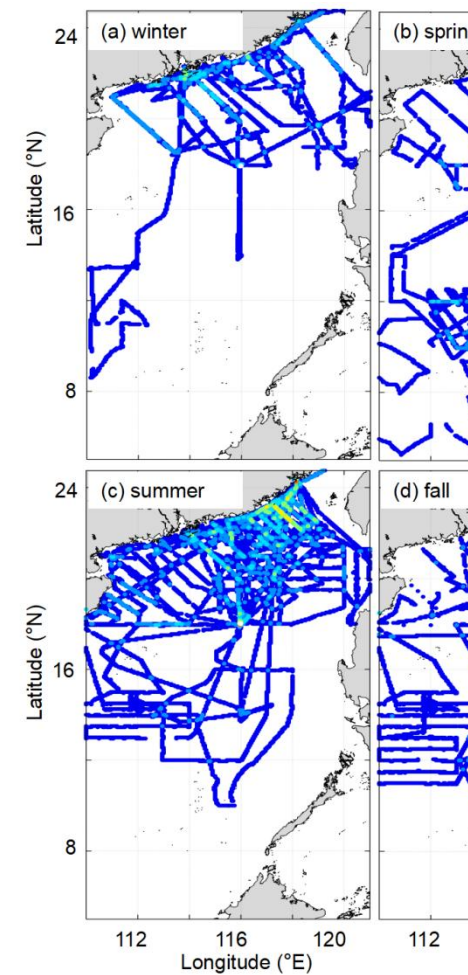
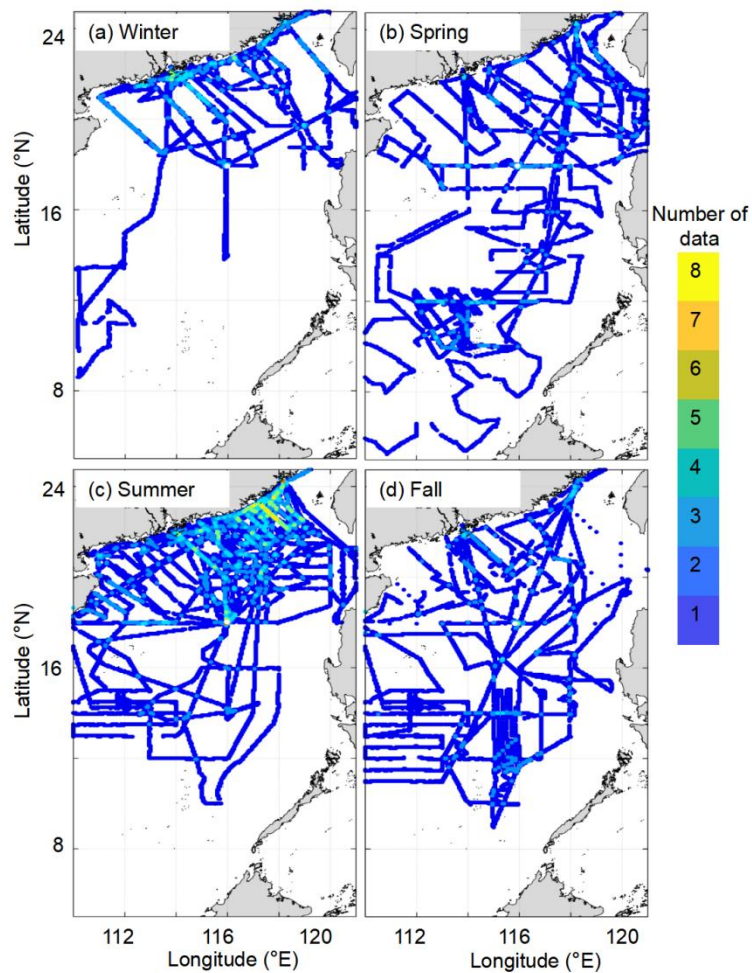


Figure 2. Cruise tracks of the observations conducted in the South China Sea in each season from 2000 to 2020: (a) Winter, (b) Spring, (c) Summer, and (d) Fall. The data collected before February 2018 are from Li et al. (2020), except those collected in July 2015 and June 2017.

Figure 3 shows the spatial and temporal distributions of in situ sea surface  $p\text{CO}_2$ . Seasonally, the lowest  $p\text{CO}_2$  occurs in January, and the highest concentrations occur in May and June. Spatially, the  $p\text{CO}_2$  distribution in the basin is relatively homogeneous, although is highly variable in the northern region. In the northern coastal area in summer, the  $p\text{CO}_2$  distribution is affected by the Pearl River plume (yielding low values) and coastal upwelling (yielding high values), which last into early fall. In winter and early spring, relatively low  $p\text{CO}_2$  values ( $\sim 350 \mu\text{atm}$ ) were found in the near-shore area. In addition, the high  $p\text{CO}_2$  values recorded on the western side of the Luzon Strait in December demonstrate the influence of winter upwelling during some of the surveys.

In addition to the above in situ sea surface  $p\text{CO}_2$  data, we selected in situ sea surface  $p\text{CO}_2$  data collected during four independent surveys across the four seasons: September 2018 (fall), December 2018 (winter), August 2019 (summer), and April 2020 (spring) to verify the accuracy of our reconstruction model in extrapolating periods lacking training datasets. Furthermore, we used an additional dataset of sea surface  $p\text{CO}_2$  calculated from observed dissolved inorganic carbon and total alkalinity during 2003–2019

删除[Author]:

删除[Author]:

删除[Author]: w

删除[Author]: s

删除[Author]: s

删除[Author]: f

删除[Author]: were

删除[Author]: ,

删除[Author]: which are from Li et al. (2020)

删除[Author]: water

删除[Author]: of in situ measurements

设置格式[admin]: 字体: 倾斜

设置格式[admin]: 下标

删除[Rick Smith]: with large

删除[Rick Smith]: i

删除[Rick Smith]: ity

删除[Author]: Spatially, the  $p\text{CO}_2$  distribution in the basin ...

删除[Author]: observed

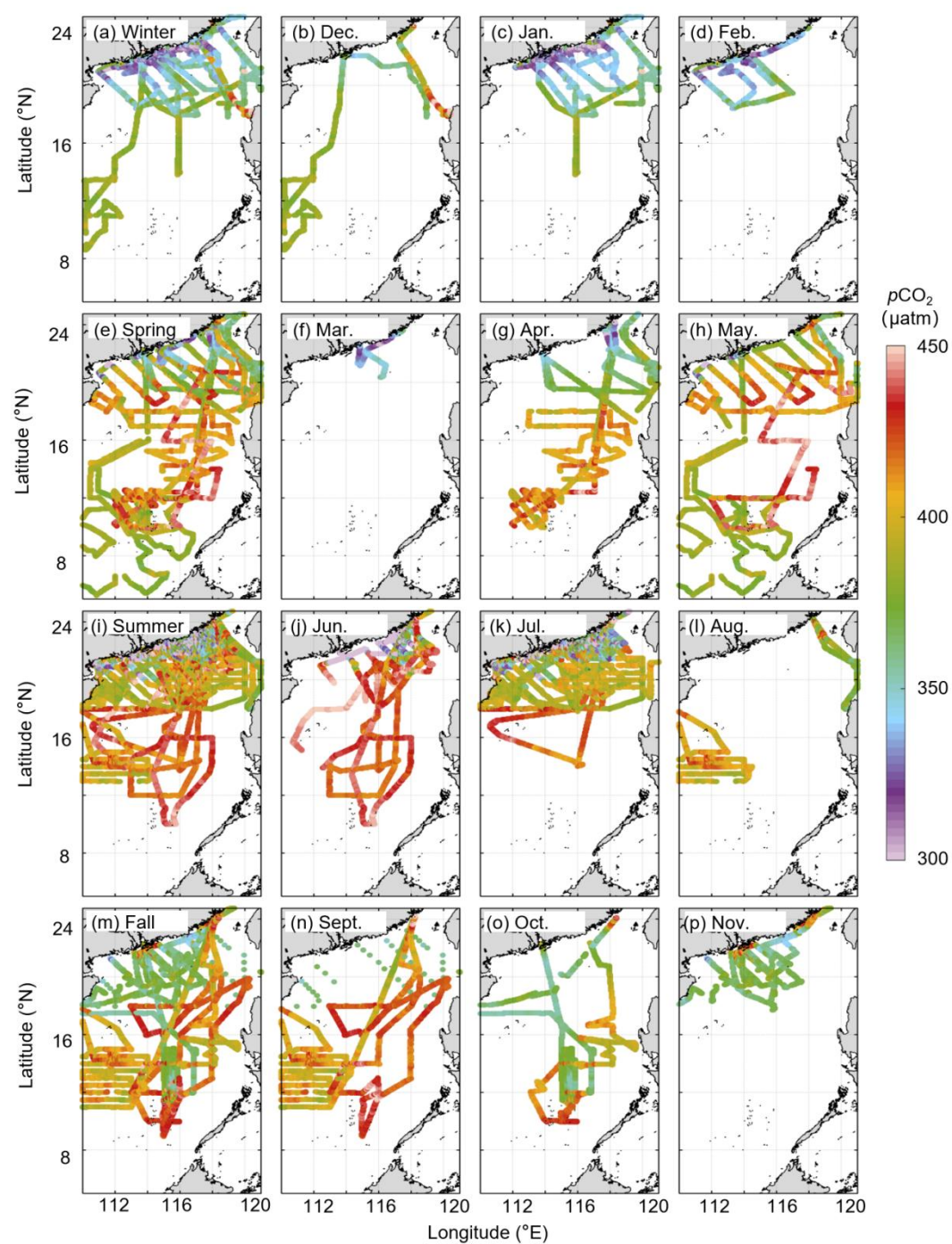
删除[admin]: in situ

删除[Author]: data



157 at the Southeast Asia Time-Series (SEATs) station (data from Dai et al., 2022) to test the long-term consistency of the

158 reconstruction.



159

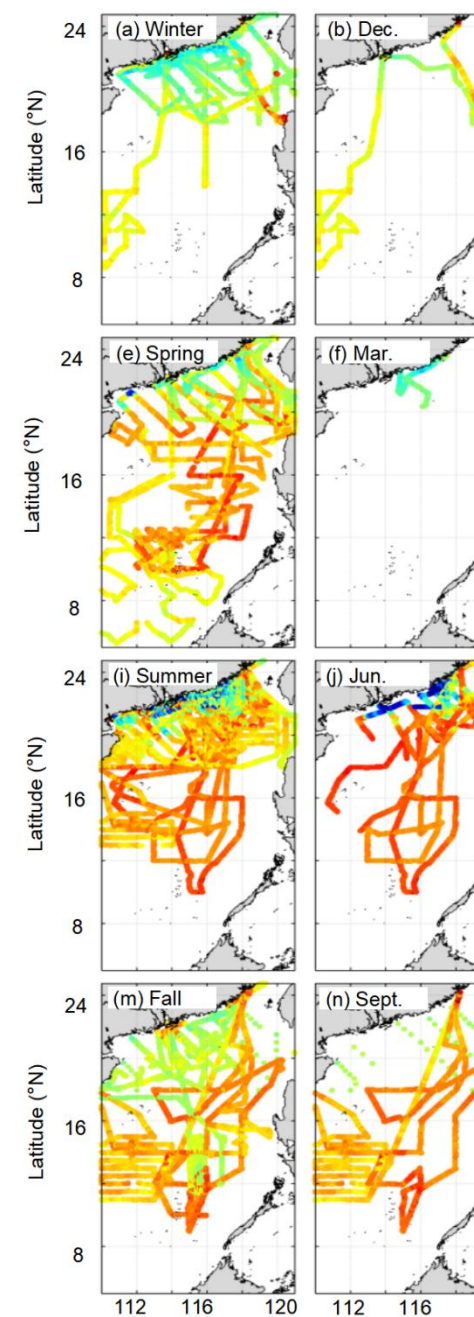
160 Figure 3. Seasonal and monthly sea surface  $p\text{CO}_2$  fields in the South China Sea; **a. Winter; b. December; c. January; d.**

161 **February; e. Spring; f. March; g. April; h. May; i. Summer; j. June; k. July; l. August; m. Fall; n. September; o. October;**

162 **p. November.** The data sources are given in Table 1.

163

164 **2.3 Remote sensing-derived sea surface  $p\text{CO}_2$  data**



删除[Author]:

删除[Author]: .

删除[Author]: The data sources can be found in Table 1 (

删除[Author]: w

删除[Author]: )

删除[Rick Smith]: can be found



The gridded (0.05°×0.05°) **RS-derived**  $p\text{CO}_2$  data cover almost the entire SCS (5–25° N, 109–122° E), and show major variations in sea surface  $p\text{CO}_2$  at the basin scale (Wang et al., 2021; Yu et al., 2022). Further details of the **RS-derived**  $p\text{CO}_2$  data can be found on the SatCO<sub>2</sub> platform (www.SatCO2.com).

A grid-to-grid comparison was undertaken between the **RS-derived**  $p\text{CO}_2$  and the **in situ**  $p\text{CO}_2$  data (Table 2). The differences in between range from 35 to 120  $\mu\text{atm}$  in the **near-shore** area. The largest biases occur in summer, when the **RMSE** is up to 29.95  $\mu\text{atm}$  (Table 2). Relatively large discrepancies may reflect the limitations of the current algorithm (MeSAA and non-linear regression), which only considers biological processes and the turbidity induced by the Pearl River discharge (characterized by Chl *a* and the remote sensing reflectance at 555 nm ( $r_{rs555}$ ), and does not take into account the riverine dissolved inorganic carbon and the input of other substances that may affect  $p\text{CO}_2$  (Bai et al., 2015, Yu et al., 2022 and Wang et al., 2021)).

To remove the influence of the bias in **RS-derived**  $p\text{CO}_2$  data on our reconstructed results, this study used the EOF method to compute the spatial patterns of the **RS-derived**  $p\text{CO}_2$  data as input data instead of directly using the **RS-derived**  $p\text{CO}_2$  data. Moreover, using EOF modes of the **RS-derived**  $p\text{CO}_2$  as input data in the reconstructed model can provide spatial constraints on the  $p\text{CO}_2$  reconstruction.

**Table 2. Biases between the seasonal remote sensing-derived  $p\text{CO}_2$  data and in situ  $p\text{CO}_2$  data, and between the reconstructed and the in situ  $p\text{CO}_2$  data. (unit:  $\mu\text{atm}$ ; the remote sensing-derived  $p\text{CO}_2$  data during 2003–2019 are from www.SatCO2.com and the source of in situ data can be found in Table 1. The reconstructed  $p\text{CO}_2$  data are from section 3; all data were gridded into 0.05°×0.05°; / means no data). MAE = mean absolute error; RMSE = root mean square error; R<sup>2</sup> = coefficient of determination; MAPE = mean absolute percentage error.**

		<b>RS-derived</b> <b><math>p\text{CO}_2</math> data</b>	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 <sup>2</sup>	/	0.98	0.92	/	/
	MAPE	/	0.01	0.01	/	/
Summer	MAE	16.75	2.48	8.46	5.73	/
	RMSE	29.95	3.54	14.69	15.18	/
	R <sup>2</sup>	/	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 <sup>2</sup>	/	0.98	0.92	/	/

删除[admin]: remote sensing

删除[admin]: ed

删除[Author]: the

删除[Author]: CO<sub>2</sub>

删除[Author]: at a

删除[Rick Smith]: on

删除[Author]: large

删除[Author]: Bai

删除[Author]: unpublished

删除[Author]: In the retrieval algorithm of Yu et al. (202 ...)

删除[Author]:

删除[Author]: remote sensing (RS) data

删除[Author]: i

删除[Author]: (Fig. 4)

删除[Author]: and the RMSE of

删除[Author]:

删除[Author]: dataRS data-derived  $p\text{CO}_2$

删除[Author]: values were compared with

删除[Author]: observed

删除[Author]: data

删除[admin]: This comparison shows that the

删除[Rick Smith]: between the RS- derived  $p\text{CO}_2$  dataRS ...)

删除[Author]: data

删除[Rick Smith]: s

删除[Author]: coastal

删除[admin]: , and that t

删除[Rick Smith]: . TIn terms of the RMSE (Table 2), the ...)

删除[Author]: bias

删除[Author]: reaches

删除[Author]: 30.0

删除[Rick Smith]: in summer

删除[Author]: Bai et al. (2015), Yu et al. (2022); unpublis ...)

删除[Author]: pointed out that r

删除[admin]: only

	MAPE	/	0.01	0.01	/	/
	MAE	9.25	2.18	5.61	11.41	/
Winter	RMSE	14.26	3.14	8.82	12.63	/
	R <sup>2</sup>	/	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 <sup>2</sup>	/	0.99	0.91	/	/
	MAPE	/	0.01	0.01	/	/

## 2.4 Other data

The RS-derived SST data produced by MODIS (<https://oceancolor.gsfc.nasa.gov/>) are adopted in our reconstruction. The uncertainty of this dataset in the SCS is  $\sim 0.27^\circ$  (Qin et al., 2014). For sea surface salinity (SSS) data, Wang et al. (2022) found relatively large differences between different open source SSS databases (i.e., multi-satellite fusion data from <https://podaac.jpl.nasa.gov/>; model data from <https://climatedataguide.ucar.edu/>; multidimensional covariance model data from <https://resources.marine.copernicus.eu/>) and the in situ SSS data. Thus, Wang et al. (2022) produced an RS-derived SSS database using machine learning methods based on the MODIS-Aqua remote sensing data. The bias between the RS-derived SSS (Wang et al., 2022) and in situ data was near-zero (mean absolute error, MAE:  $\sim 0.25$ ). Next, we used Chl-*a* (from <https://oceancolor.gsfc.nasa.gov/>) as an indicator of biological influence, which has a bias of  $\sim 0.35$  on a log scale and  $\sim 115\%$  in the SCS (Zhang et al., 2006). Atmospheric  $p\text{CO}_2$  also influences sea surface  $p\text{CO}_2$  through air-sea  $\text{CO}_2$  exchange. We chose the atmospheric  $\text{CO}_2$  mole fraction ( $x\text{CO}_2$ ) data from the monthly mean  $\text{CO}_2$  concentrations measured at the Mauna Loa Observatory, Hawaii (<https://gml.noaa.gov/>), and then calculated the atmospheric  $p\text{CO}_2$  values from  $x\text{CO}_2$  using the method of Li et al. (2020).

## 3 Methods

The  $p\text{CO}_2$  reconstruction procedure is shown in Figure 4. It includes: (1) data processing and (2) model training and testing. For the former, we firstly gridded the in situ data and RS-derived  $p\text{CO}_2$  data into  $0.05^\circ \times 0.05^\circ$  boxes with a monthly temporal resolution. Secondly, we filled missing  $p\text{CO}_2$  measurements with the RS-derived  $p\text{CO}_2$  data according to Fay et al. (2021) (see more details in Section 3.1). We then used EOF to ignore any biases in the RS-derived  $p\text{CO}_2$  dataset itself or from the  $p\text{CO}_2$  filling method. Thirdly, the gridded in situ  $p\text{CO}_2$  data and their corresponding RS-derived data were divided into a training set (90%) and a testing set (10%) to calculate the  $p\text{CO}_2$  retrieval model. To ensure that the model had sufficient training samples in the coastal area, we divided the entire SCS into two regions along the 200 m isobath (as shown in Figure 5). The data from these two regions

删除[Author]:

删除[Author]:

删除[Author]: here

删除[Author]: the

删除[Author]: SSS

删除[Author]: .

删除[Author]: the

删除[Author]: .Wang et al. (in preparation2022) found a

删除[Author]: in preparation

删除[Author]: reconstructedproduced

删除[Author]: remote sensing

删除[Author]:

删除[Author]: by

删除[Author]: based on based on a combination of

删除[Author]: remote sensing

删除[Author]: reconstructed

删除[Author]: observed

删除[Author]: . Chl-*a* data from MODIS

删除[Rick Smith]: have

删除[Author]: water

删除[Rick Smith]: atmosphere

删除[Author]: .

删除[Author]: T

删除[Author]: were calculated

删除[Author]: by

删除[Author]: 5

删除[Author]: observed

删除[Author]:

删除[Author]: RS  $p\text{CO}_2$  data

删除[Rick Smith]: grid

删除[Author]: And all these data used in machine learning

删除[Author]: used the  $p\text{CO}_2$  filling method according to

删除[Author]: the

删除[Author]:

205

were divided into training and testing sets with the same ratios listed above (9:1), and then combined to obtain the final training

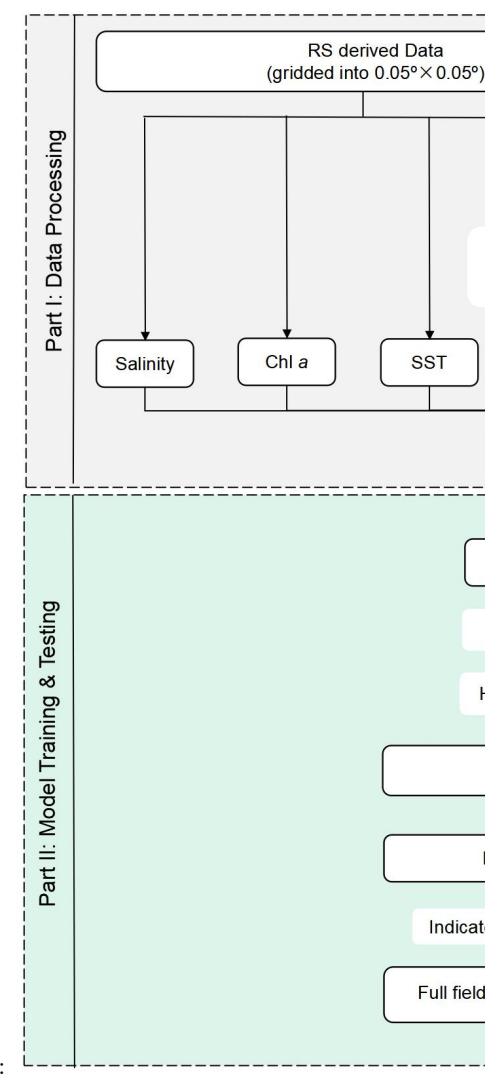
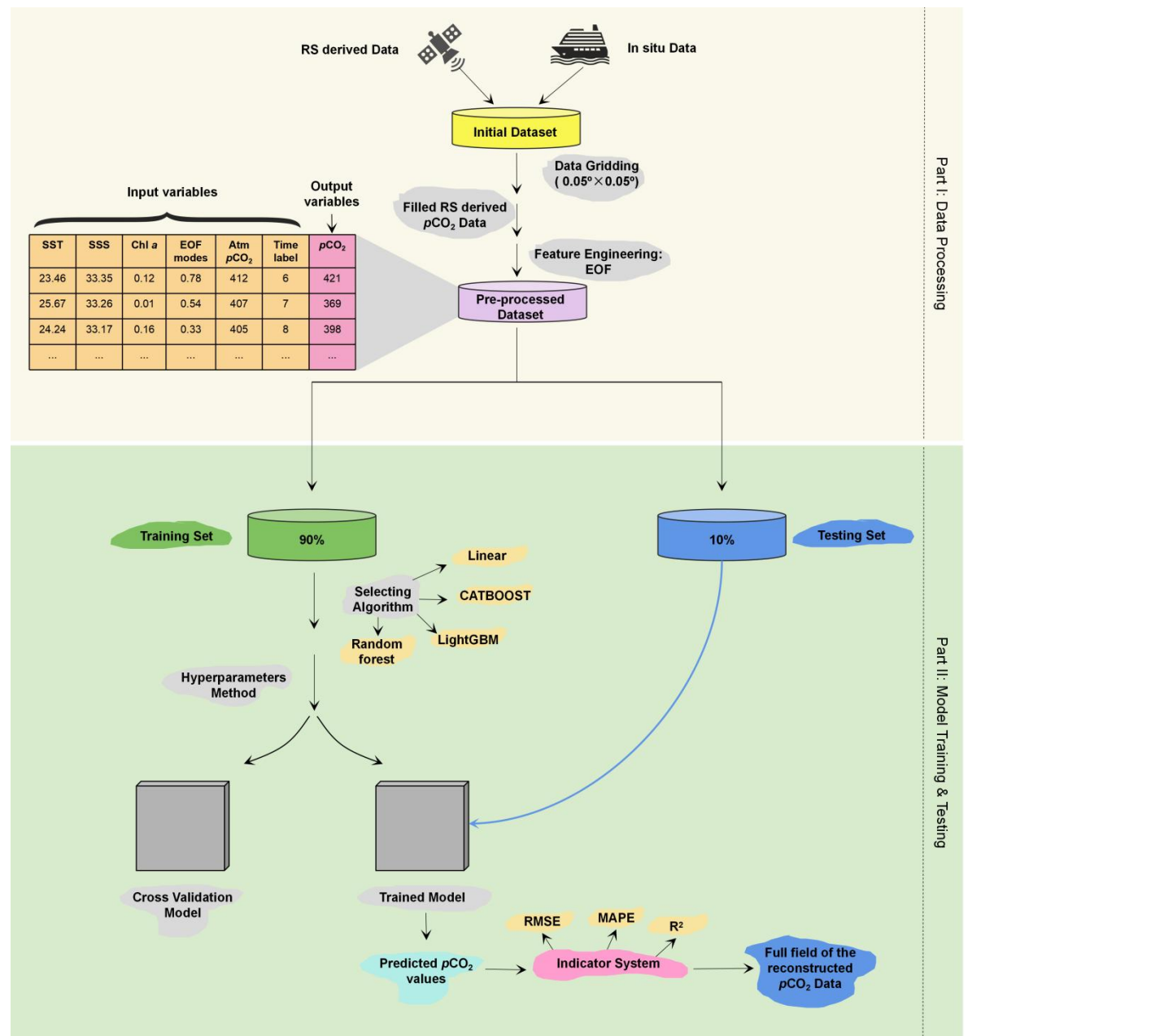
删除[Author]: which were

206

and testing sets. Note that all the data used in the machine learning have been interpolated on the same grid.

删除[Author]: And

删除[Rick Smith]: se



删除[Author]:

207

Figure 4. Procedure for the reconstruction of surface water pCO<sub>2</sub> using machine learning. RS-derived data = remote sensing derived data, RMSE = root mean square error, MAPE= mean absolute percentage error, and R<sup>2</sup> = coefficient of determination, and MAE = mean absolute error.

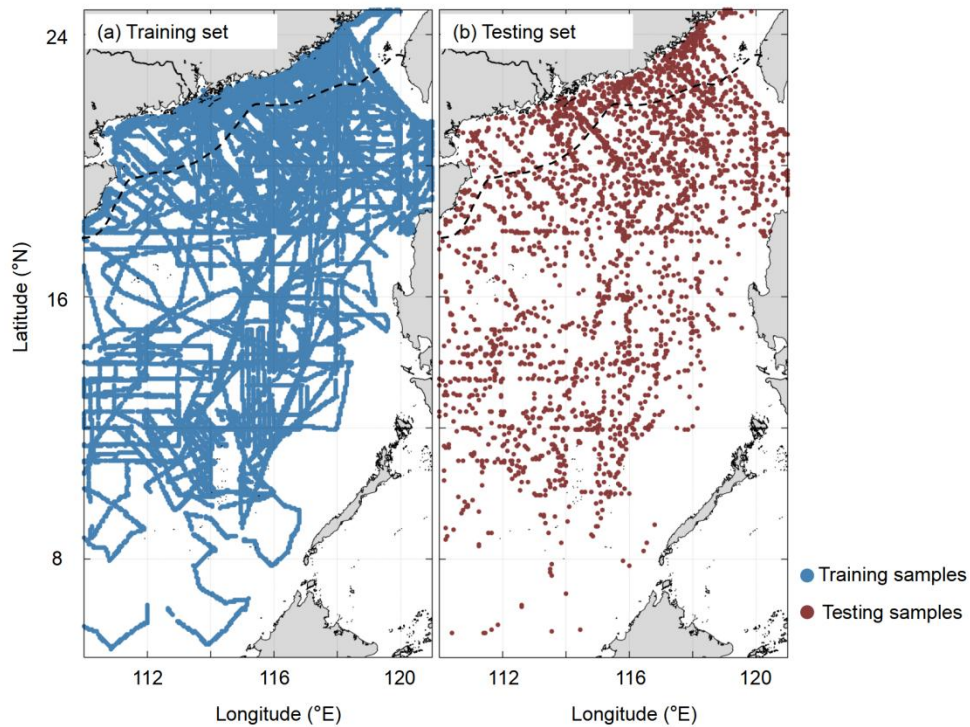
删除[Author]: 5

删除[Author]:

删除[Author]:

删除[Rick Smith]: average





**Figure 5. Spatial distributions of Training samples (a) and Testing samples (b); The black dashed line shows the 200 m isobath.**

For model training and testing, we chose a relatively reliable algorithm to undertake the  $p\text{CO}_2$  reconstruction. Next, we determined the optimal range of the parameters using hyperparameter methods (code from <https://github.com/optuna/>) for the training set. The final optimal parameter values were then determined using the  $K$ -fold and cross validation method (code from <https://github.com/suryanktiwari/Linear-Regression-and-K-fold-cross-validation>) for the training set. These optimal parameters were applied to the chosen algorithm. Finally, the testing set was used to verify the accuracy of the  $p\text{CO}_2$  retrieval algorithm produced by the training set, and some indicators of the model's accuracy were calculated. More detailed methods employed in the present study are described below.

### 3.1 Remote sensing data filling

As mentioned in the SatCO2 platform ([www.SatCO2.com](http://www.SatCO2.com)), RS-derived  $p\text{CO}_2$  datasets have some missing values. Thus, we used the  $p\text{CO}_2$  data filling method suggested by Fay et al. (2021) to obtain the missing datapoints. First, a scaling factor for a filled month was calculated according to Equation 1:

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

where  $sf_{pCO_2}$  is the scaling factor,  $pCO_2^{ens}$  is the monthly RS-derived  $p\text{CO}_2$  data, and  $pCO_2^{lim}$  is the monthly climatology RS-derived  $p\text{CO}_2$  data;  $x$  and  $y$  indicate that we took the area-weighted average over longitude ( $x$ ) and latitude ( $y$ ) to produce the monthly  $sf_{pCO_2}$  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. (2021) for details of this method).

删除[Author]: The s

删除[Author]: o

删除[Rick Smith]: stands for

删除[Rick Smith]: depth contour.

删除[Rick Smith]: firstly

删除[Rick Smith]: After that

删除[Rick Smith]: e

删除[Rick Smith]: the

删除[Author]:

删除[Zhixuan Wang]: RS  $p\text{CO}_2$  data may

删除[Author]: are missing

删除[Rick Smith]: misses some

删除[Rick Smith]: fill in

删除[Rick Smith]: portions

删除[Author]:

删除[Author]: RS  $p\text{CO}_2$  datum

删除[Author]:

删除[Author]: RS  $p\text{CO}_2$  datum

230 Briefly, this filling method scales the climatological monthly  $p\text{CO}_2$  field values to fill in the missing measurements. Therefore,  
231 although specific values may be biased, the interpolated measurements still retain the main spatial distribution pattern of the filled  
232 months.

### 233 3.2 Feature engineering and selection

234 As mentioned above, the  $p\text{CO}_2$  data filling method may bias some of the actual values. To avoid the influence of such biases on the  
235 reconstructed results, instead of directly using the RS-derived  $p\text{CO}_2$  data as features in our reconstructed model, we used the EOF  
236 method to obtain the main spatiotemporal distribution patterns of the RS-derived  $p\text{CO}_2$  data as features in our reconstructed model.  
237 The EOF reflects the spatial commonality of variables shown in the time series, and thus it is widely used to calculate spatial  
238 patterns of climate variability (e.g. Levitus et al., 2005; Dye et al., 2020; McMonigal and Larson, 2022). Typically, the spatial  
239 commonality of variables (EOF modes) is found by computing the eigenvalues and eigenvectors of a spatially weighted anomaly  
240 covariance matrix of a field. Each EOF modes' corresponding variance represents its degree of interpretation of the spatial pattern  
241 of a variable. For each of the 12 months, the cumulative variance contribution of the first eight EOF values was consistently >  
242 90%, indicating that it could explain the main  $p\text{CO}_2$  spatial characteristics during each month, we therefore selected them as  
243 features.

244 The features selected in our reconstructed model can be divided into two main categories. In the first category, the features are  
245 related to the underlying physicochemical mechanisms controlling the  $p\text{CO}_2$  distribution: for example, that SST exerts a primary  
246 control on the seasonal variations in surface water  $p\text{CO}_2$  in the northern SCS (Zhai et al., 2005; Chen et al., 2007; Li et al., 2020).  
247 In the second category, they provide spatiotemporal information for the  $p\text{CO}_2$  reconstruction. Previous studies (Landschützer et al.,  
248 2014; Laruelle et al., 2017; Denvil et al., 2019) have shown that Chl-*a* plays a critical role in fitting the influence of biological  
249 activity to  $p\text{CO}_2$ , especially in the northern SCS (Landschützer et al., 2014; Laruelle et al., 2017; Denvil et al., 2019). Sutton et al.  
250 (2017) suggest that increasing atmospheric  $p\text{CO}_2$  controls the overall increase in seawater  $p\text{CO}_2$ . For the features that provide  
251 spatiotemporal information for the  $p\text{CO}_2$  reconstruction, in the present study we selected the first eight EOF values of  $p\text{CO}_2$  as the  
252 main spatial distribution feature and monthly information of the in situ datasets as the temporal feature.

### 253 3.3 Algorithm selection

254 Ensemble learning, which is the process of training multiple machine learning models and combining their output to improve the  
255 reliability and accuracy of predictions, is one of the most powerful machine learning techniques (e.g., Zhan et al., 2022; Chen et  
256 al., 2020). (e.g., Zhan et al., 2022; Chen et al., 2020). In other words, several different models are used as the basis to develop an  
257 optimal predictive model. There are two main ways to employ ensemble learning: bagging (to decrease the model's variance), or  
258 boosting (to decrease the model's bias). The random forest algorithm (code from <https://scikit-learn.org/stable/>) is an extension of  
259 the bagging method as it utilizes both bagging and feature randomness to create an uncorrelated forest of decision trees. The Light  
260 Gradient Boosting Machine (LightGBM; code from <https://github.com/microsoft/LightGBM/>) is a gradient boosting framework

设置格式[Rick Smith]: 非上标/ 下标

删除[Author]:

删除[Author]: RS  $p\text{CO}_2$  data

删除[Author]: feengineered featured data (via the

删除[Author]: )

删除[Author]:

删除[Author]: RS  $p\text{CO}_2$  data

删除[Rick Smith]:

删除[Rick Smith]: , also named

删除[Rick Smith]:

删除[Rick Smith]: ,

删除[Author]: are

删除[Author]: The EOF reflects the spatial commonality ...

删除[Rick Smith]: it

删除[Rick Smith]: ,

删除[Rick Smith]: and

删除[Rick Smith]: selection

删除[Rick Smith]: T

删除[Rick Smith]: one

删除[Rick Smith]: is

删除[Rick Smith]: , and

删除[Rick Smith]: other one

删除[Rick Smith]: can

删除[Zhixuan Wang]: . For example, that SST exerts a ...

删除[Rick Smith]: researches

删除[Author]: s

删除[Rick Smith]: the

删除[Rick Smith]: e in

删除[Rick Smith]: whereas

删除[Author]: observed

删除[Rick Smith]: provides

删除[Rick Smith]: It is the process of training multiple ...

删除[Rick Smith]: Different

that uses tree-based learning algorithms. LightGBM can be used for regression, classification, and other machine learning tasks; it exhibits rapid, high-performance as a machine learning algorithm. CATBOOST (code from <https://github.com/catboost/>) is a gradient boosting algorithm, which improves prediction accuracy by adjusting weights according to the data distribution and by incorporating prior knowledge about the dataset. This can help to reduce overfitting and improve general performance.

From the above options, we chose three ensemble learning algorithms as the machine learning-based regression portion, and multi-linear regression methods (Wang et al., 2021) as the linear regression portion. We then used the K-fold and cross validation methods to verify the applicability of different regression algorithms in the  $p\text{CO}_2$  reconstruction for seasonal training data. The results show that in summer, the CATBOOST algorithm yields the best degree of accuracy, with an RMSE of  $16 \mu\text{atm}$  (Table R1). In contrast, the RMSE of LightGBM was  $27 \mu\text{atm}$ , and that of Random Forest was  $26 \mu\text{atm}$ . The RMSE was nearly  $20 \mu\text{atm}$  using the linear regression algorithm employed by Wang et al. (2021). Thus, CATBOOST appears to provide a reliable algorithm for reconstructing  $p\text{CO}_2$ . In the other three seasons, however, using different algorithms resulted in minor differences ( $\sim 2 \mu\text{atm}$  in RMSE).

**Table 3. RMSEs associated with different algorithms in the four seasons.**

Season	Random Forest	LightGBM	CATBOOST	Multi-linear regression (Wang et al., 2021)
Spring	$10.65 \mu\text{atm}$	$9.52 \mu\text{atm}$	$8.17 \mu\text{atm}$	NaN*
Summer	$26.53 \mu\text{atm}$	$27.83 \mu\text{atm}$	$16.15 \mu\text{atm}$	$20.13 \mu\text{atm}$
Fall	$10.34 \mu\text{atm}$	$11.56 \mu\text{atm}$	$10.35 \mu\text{atm}$	NaN
Winter	$12.48 \mu\text{atm}$	$12.75 \mu\text{atm}$	$11.52 \mu\text{atm}$	NaN

\*NaN stands for missing values

### 3.4 Evaluation metrics

It is necessary to evaluate the accuracy of any model based on certain error metrics before applying it to specific scenarios.

Common model evaluation metrics include RMSE, MAPE,  $R^2$  (coefficient of determination), and MAE.

The mean squared error (MSE) is the standard deviation of the residuals (prediction error), and the residuals are the distances between the fitted line and the data points (i.e., the residuals show the degree of concentration of the reconstructed data around the regression line. In regression analysis, RMSE is commonly used to verify experimental results. To assess bias, the RMSE needs to combine the magnitude of the model data and is calculated as:

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

where  $y$  stands for the in situ data,  $y_r$  represents the reconstructed data, and  $n$  is the number of datapoints.

The mean absolute percentage error (MAPE) is a statistical measure used to define the accuracy of a machine learning algorithm on a particular dataset. It is commonly used because, compared to other metrics, it uses a percentage to measure the magnitude of

删除[Rick Smith]: and

删除[Zhixuan Wang]: ization

删除[Author]: the

删除[Author]: We

删除[Author]: ,

删除[Author]: ,

删除[Author]: For comparison

删除[Author]:

删除[Author]: Note that

删除[Author]: for other three seasons only

删除[Author]: From the above options, we chose three ensemble learning algorithms as the machine learning-based regression portion, and multi-linear regression methods (Wang et al., 2021) as the linear regression portion, and we then used the K-fold and cross validation methods to verify the

删除[Author]: between

删除[Rick Smith]: different

删除[Author]: The RMSE between the CO2 different algorithm and in situ data of different seasonal (NaN stand for the missing value

删除[Rick Smith]: the

删除[Author]:

删除[Rick Smith]: stands for

删除[Rick Smith]: where

删除[Rick Smith]: represent

删除[Rick Smith]: ,

删除[Author]: .

删除[Author]: it

删除[Rick Smith]: stands for

删除[Author]: observational data



the bias and is easy to understand and interpret; the lower the value of the MAPE, the better a model is at forecasting. MAPE is calculated as follows:

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

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

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

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

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

### 3.5 Uncertainty

In previous studies, RMSE and MAE have primarily been used to represent the uncertainties in reconstructed datasets. However, this expression of uncertainty ignores the sensitivity of the reconstructed model to the features; i.e., the biases that the features themselves pass to the reconstructed model are ignored. Moreover, it is clearly unreasonable to use a single RMSE or MAE value to represent the entire region because the spatial bias pattern in the coastal region clearly differs from that in the basin. Thus, here we present a novel method for calculating uncertainty, as shown below:

$$\text{Uncertainty} = \text{MAX} \left( \left[ \frac{\sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \frac{|OR\_Monthly\_Data(i,j,k) - Obs\_Monthly\_Data(i,j,k)|}{Obs\_Monthly\_Data(i,j,k)}}{\text{num}(i) + \text{num}(j)}, \dots, \frac{\sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \frac{|OR\_Monthly\_Data(i,j,k) - Obs\_Monthly\_Data(i,j,k)|}{Obs\_Monthly\_Data(i,j,k)}}{\text{num}(i) + \text{num}(j)} \right] \right) * 100\% * pCO_2 \text{ recon} + \left( \frac{\partial pCO_2}{\partial Feature} \right) dFeature \quad (6)$$

Equation (6) includes two terms: the first term is the conservative bias between the reconstructed  $pCO_2$  fields and the in situ data, and the second is the sum over sensitivity of the reconstructed model to the features. For the first term in Equation 6,  $k$  stands for the  $k$ th month,  $OR\_Monthly\_Data(i, j, k)$  stands for the  $k$ th monthly reconstructed data at longitude ( $i$ ) and latitude ( $j$ ), and  $Obs\_Monthly\_Data(i, j, k)$  stands for the  $k$ th monthly in situ data at longitude ( $i$ ) and latitude ( $j$ ). Therefore,  $MAX$  in the first term stands for the maximum of the  $k$  monthly bias ratios. And ' $pCO_2 \text{ recon}$ ' stands for the reconstructed  $pCO_2$  data. In the second term, where  $dFeature$  stands for the bias of the features. We conducted a sensitivity analysis using a chain rule to evaluate the influence of these biases in the features on  $pCO_2$ . Then we estimated  $pCO_2$  changes due to these features' variabilities by constraining these features based on our model, and computed  $\frac{\partial pCO_2}{\partial Feature}$ . For example, for  $\frac{\partial pCO_2}{\partial SST}$ , we only changed the value of SST

删除[Rick Smith]: ,

删除[Rick Smith]: , and

删除[Author]: field observations

删除[Author]:

删除[Rick Smith]: were

删除[Rick Smith]: mostly

删除[Author]: the

删除[Rick Smith]: ;

删除[Author]: s

删除[Author]: between

删除[Author]: coastal

删除[Author]: and basin areas

删除[Rick Smith]: we

删除[Rick Smith]: of

删除[Rick Smith]: calculation

删除[Author]:  $\langle \mathbf{A} \rangle + \langle \mathbf{B} \rangle$  (part 1)

删除[Author]: parts

删除[Author]: ;

删除[Zhixuan Wang]: (part 1 the first term)

删除[Author]: the

删除[Zhixuan Wang]: (the second term)

删除[Author]: (part 2)

删除[Author]: F

删除[Author]: R1

删除[Author]: the

删除[Author]: that

删除[Author]: value between

删除[Author]: or part 1,  $\langle \mathbf{A} \rangle$  stands for the monthly

删除[Author]:

删除[Author]: For part

删除[Author]: 2

删除[Author]: ,

删除[Author]: w

313 and kept the values of the other features constant to calculate the effect of each additional unit of SST on the simulated  $p\text{CO}_2$ .

## 314 **4 Results and discussion**

### 315 **4.1 Results**

316 The reconstructed  $p\text{CO}_2$  fields show relatively low values in the northern coastal region of the study area, and generally high

317 values in the mid and southern basins (Fig. 6). The continuous changes of the spatiotemporal distribution can be found in the

318 reconstruction results (Fig. 6). The reconstructed  $p\text{CO}_2$  fields show a trend of slow but sustained increases from 2003 to 2020.

319 Spatial patterns of  $p\text{CO}_2$  change between 2003 and 2020, such that the coastal portion of the northern SCS shows relatively

320 complex variability from multiple controlling factors, such as coastal upwelling, river plumes, biological activity, etc. However,

321  $p\text{CO}_2$  values in the mid and southern basins are relatively homogeneous, as they are mainly controlled by atmospheric  $p\text{CO}_2$

322 forcing and SST. Temporal changes in  $p\text{CO}_2$  between 2003 and 2020, are relatively large ( $\sim 44 \mu\text{atm}$ ) in summer and relatively

323 small ( $\sim 33 \mu\text{atm}$ ) in winter.

删除[Author]: ,

删除[Author]: results

删除[Author]: of the

删除[Author]: simulation

删除[Rick Smith]: study

删除[Rick Smith]: but

删除[Rick Smith]: shows

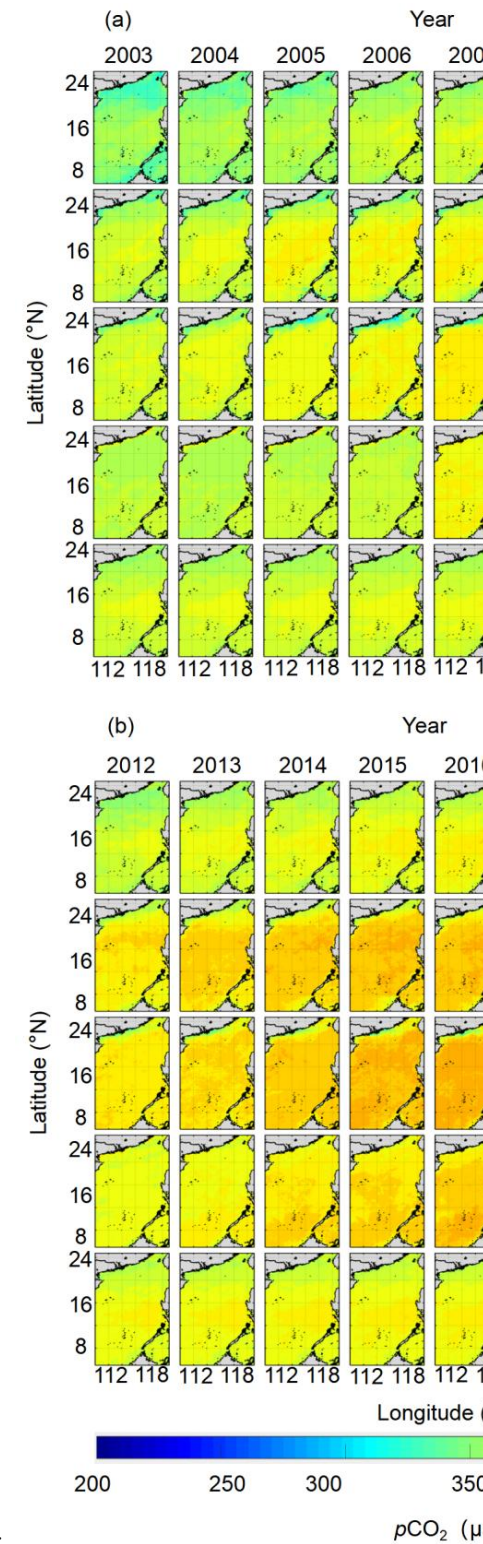
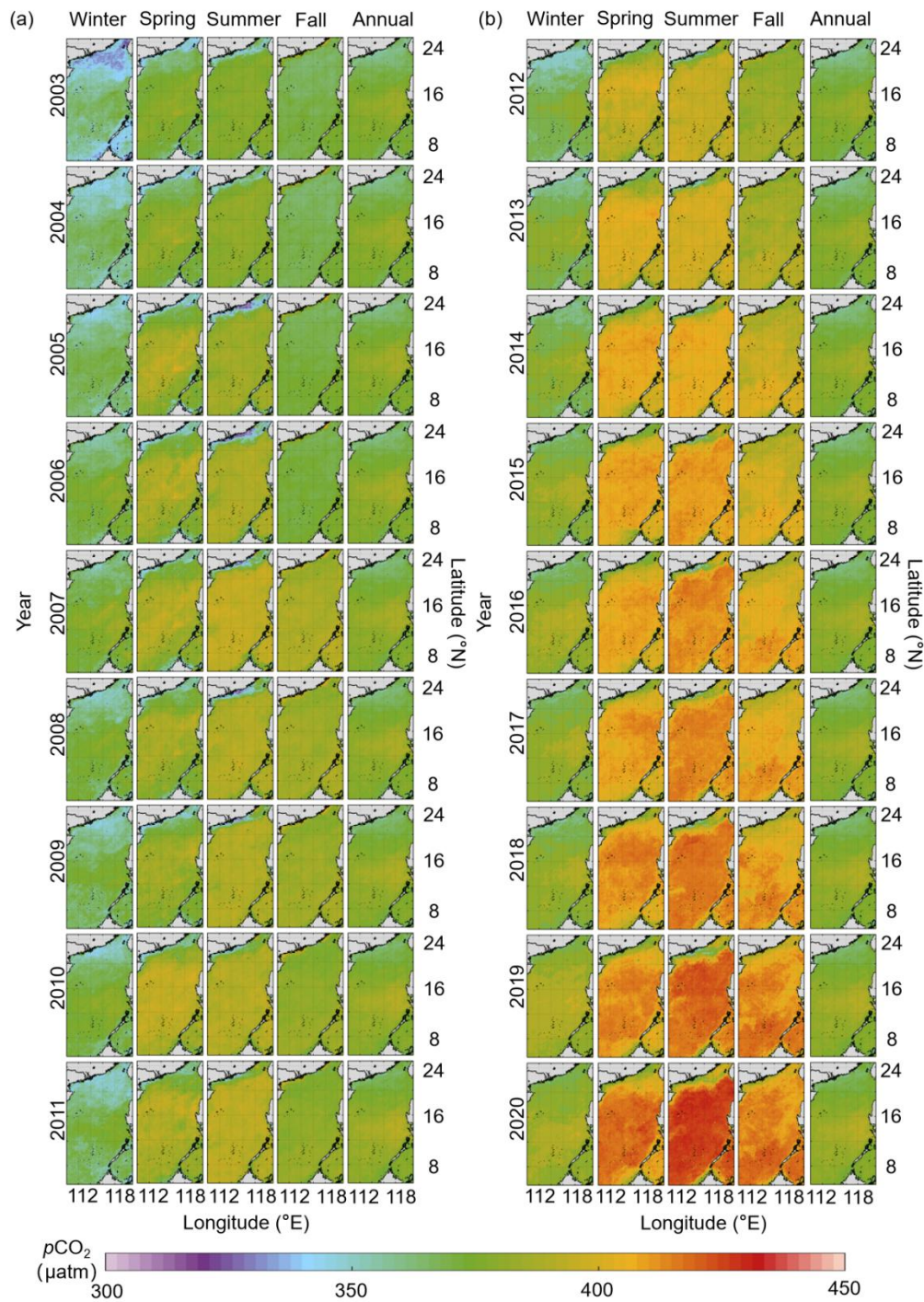
删除[Author]: continuity

删除[Rick Smith]: because

删除[Rick Smith]: of

删除[Rick Smith]: because





324  
325 **Figure 6. Reconstructed seasonal and annual  $p\text{CO}_2$  fields in the South China Sea from 2003 to 2020 (a, 2003-2011; b,**  
326 **2012-2020).**

#### 327 4.2 Model validation

328 Figure 7 compares the monthly reconstructed and in situ data. For the training dataset, the reconstructed  $p\text{CO}_2$  fields of the four  
329 seasons fit the in situ data well (Fig. 7), with an average RMSE of 3.43  $\mu\text{atm}$  and an average MAE of 2.14  $\mu\text{atm}$  (Table 2). For the  
330 testing sets, although there are some outliers, most of the reconstructed  $p\text{CO}_2$  data are consistent with the in situ data, with RMSE  
331 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 accuracies of  
332 the four seasonal models are all around 99% (Table 2), with the highest value for spring data and the lowest value for summer data.

删除[Author]:

删除[Rick Smith]: during the period

删除[Author]: field-observed

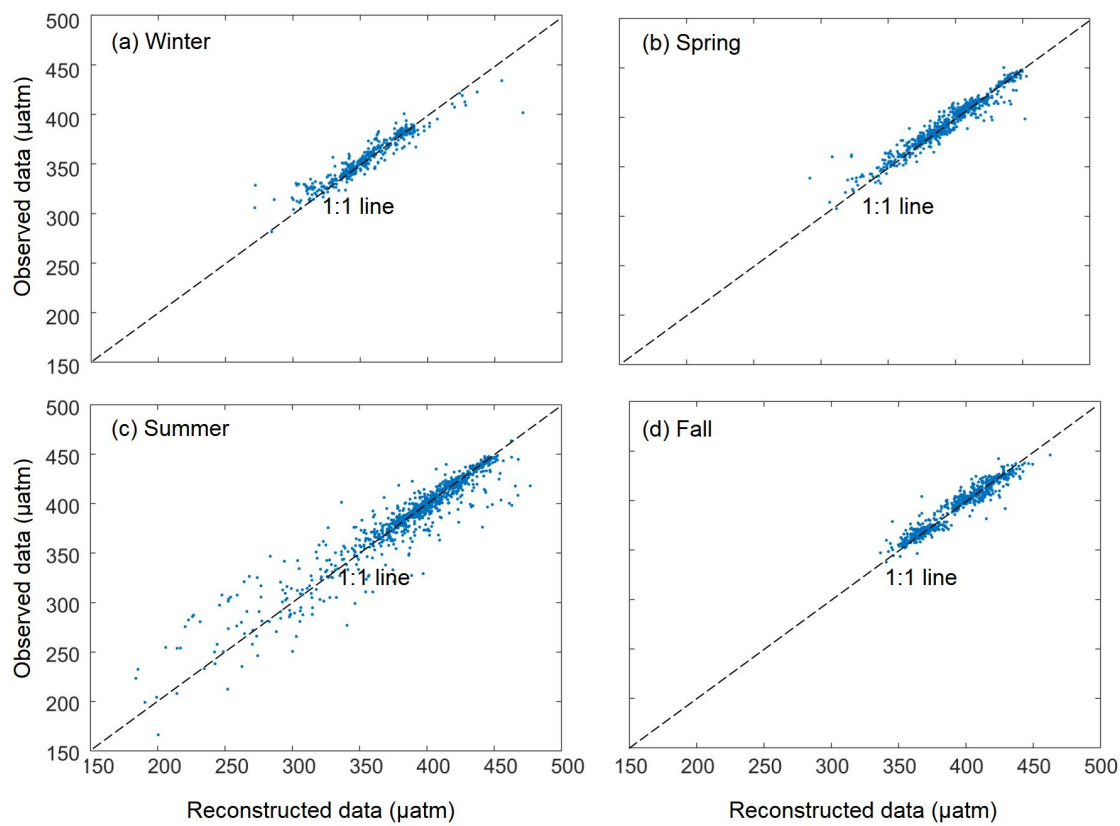
删除[Author]: field-observed

删除[Author]: ~

删除[Author]: field-observed

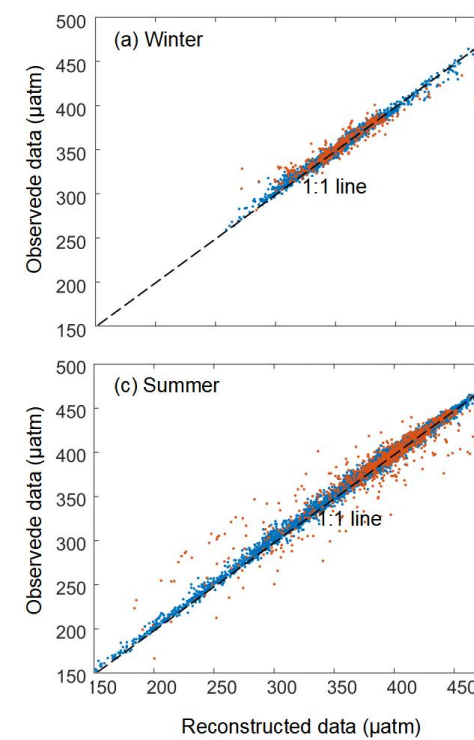


333 The relatively large bias (14.67  $\mu\text{atm}$ ) in the summer may be the influence of relatively complex regional processes, such as river  
 334 plumes and upwelling. The four evaluation metrics indicate that our reconstructed  $p\text{CO}_2$  field is highly accurate in simulating both  
 335 the training and testing sets.



删除[Author]: largestgreatest

设置格式[Zhixuan Wang]: 行距: 单倍行距



删除[Author]:

删除[Rick Smith]: the

删除[Author]: observed

删除[Author]: Ttesting

删除[Rick Smith]: were

删除[Rick Smith]: overlaid

删除[Rick Smith]: pattern

删除[Author]: field observations

删除[Rick Smith]: in

删除[Rick Smith]: distribution

设置格式[Zhixuan Wang]: 字体: (中文)宋体, 图案:

设置格式[Zhixuan Wang]: 字体: (默认)Times New Ro

删除[Zhixuan Wang]: the biases are concentrated mainly

设置格式[Zhixuan Wang]: 字体: (中文)宋体, 图案:

删除[Rick Smith]: distribution

删除[Rick Smith]: .

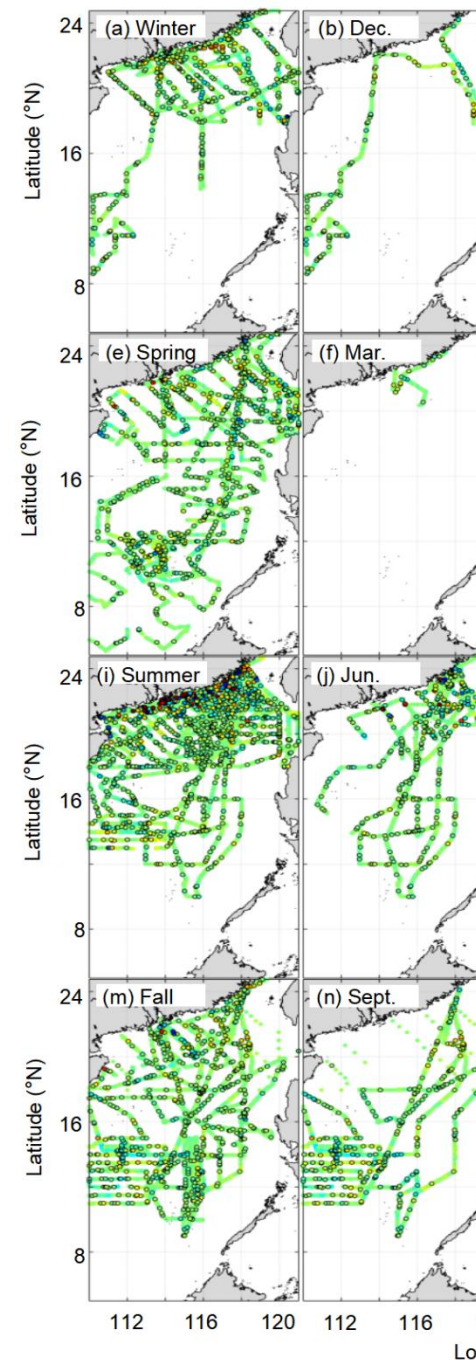
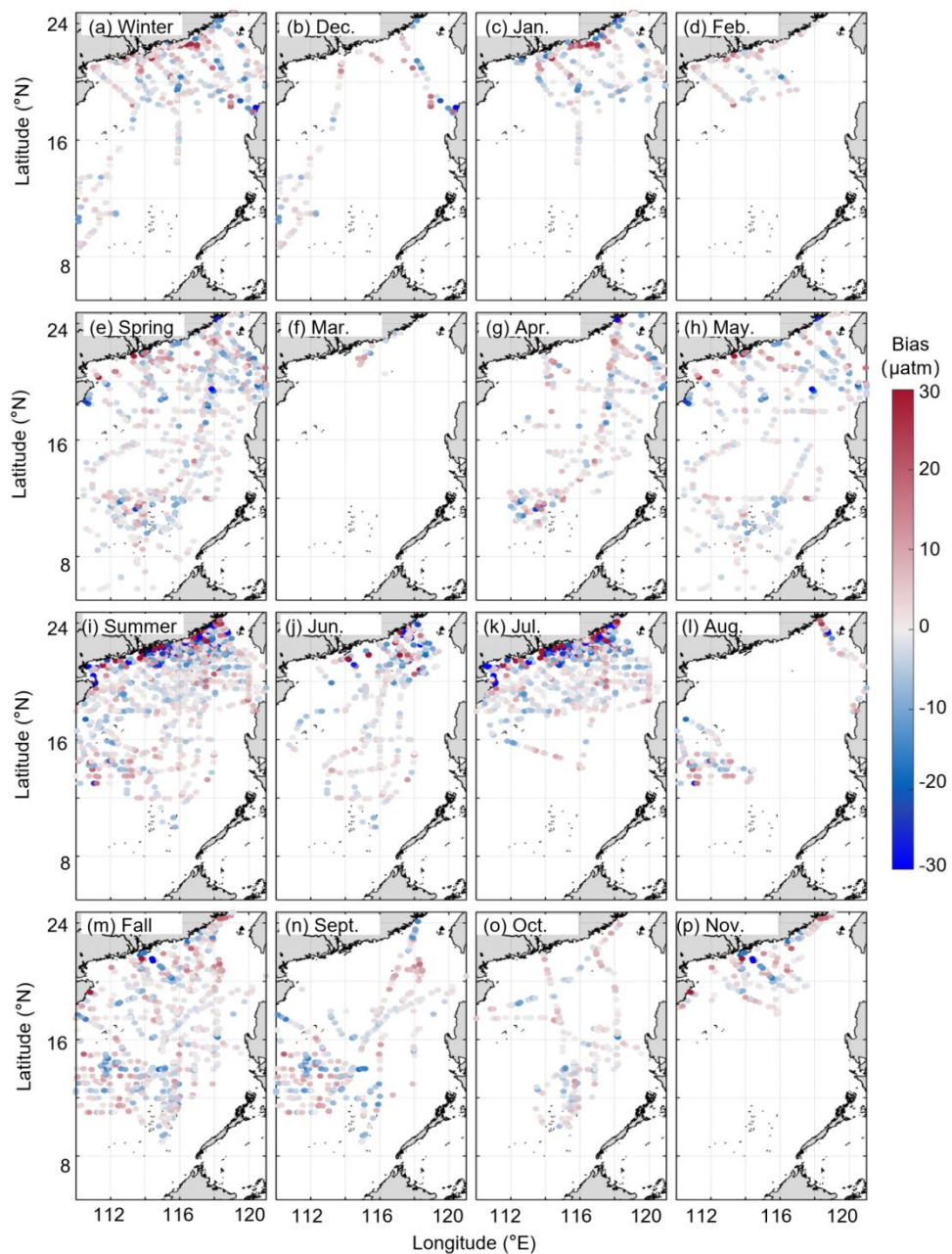
删除[Rick Smith]: T

删除[Rick Smith]: reconstruction

336  
 337 **Figure 7. Comparisons between the monthly reconstructed and in situ  $p\text{CO}_2$  values for the testing set (monthly results are  
 338 grouped into the four seasons: (a) Winter: Dec., Jan., Feb.; (b) Spring: Mar., Apr., May; (c) Summer: Jun., Jul., Aug.; (d)  
 339 Fall: Sept., Oct., Nov.).**

340 The distributions of the biases between the reconstructed fields and the in situ data for both the training and testing datasets can be  
 341 found in Figure 8. In terms of the temporal pattern, the larger biases were more concentrated in the summer. For the spatial pattern,  
 342 the biases in the northern coastal area are much greater than those in the basin. However, 95% of the biases are  $< \pm 10 \mu\text{atm}$ ;  
 343 therefore, our reconstructed dataset exhibits relatively high accuracy.

344



**Figure 8. Differences between the reconstructed and in situ  $p\text{CO}_2$  data both seasonally and monthly for the testing set (a. Winter; b. December; c. January; d. February; e. Spring; f. March; g. April; h. May; i. Summer; j. June; k. July; l. August; m. Fall; n. September; o. October; p. November).**

Figure 9 shows the bias between our reconstructed fields and the four independent *in situ* datasets corresponding to the four seasons. This validation can verify the accuracy of the *retrieval algorithm* for months without observations, namely the applicability of the *retrieval algorithm* extrapolation. This comparison shows that the *retrieval algorithm* is relatively accurate in the basin, with a near-zero bias (MAE:  $\sim 8$   $\mu\text{atm}$ , Fig. 9 a). The *largest* bias occurs in the Pearl River plume area in summer ( $\sim 35$   $\mu\text{atm}$ ). The *retrieval algorithm* also has a high accuracy for  $p\text{CO}_2$  spatial variability, except in the Pearl River plume area in summer (22–20  $^\circ\text{N}$ , Fig. 9 b–e). The effect of the Pearl River plume on the  $p\text{CO}_2$  spatial distribution in our *retrieval algorithm* is smaller than that shown by the *in situ* data. This is because at around the survey time (August 24–28, 2019), a large amount of

删除[Author]:

删除[Rick Smith]: seasonal and monthly

删除[Rick Smith]: and the observed in situ  $p\text{CO}_2$

删除[Author]: w

删除[Author]: The open circles represent the difference ...

删除[Author]: field observation

删除[Author]: e

删除[Author]: reconstruction model

删除[Rick Smith]: in

删除[Rick Smith]: data

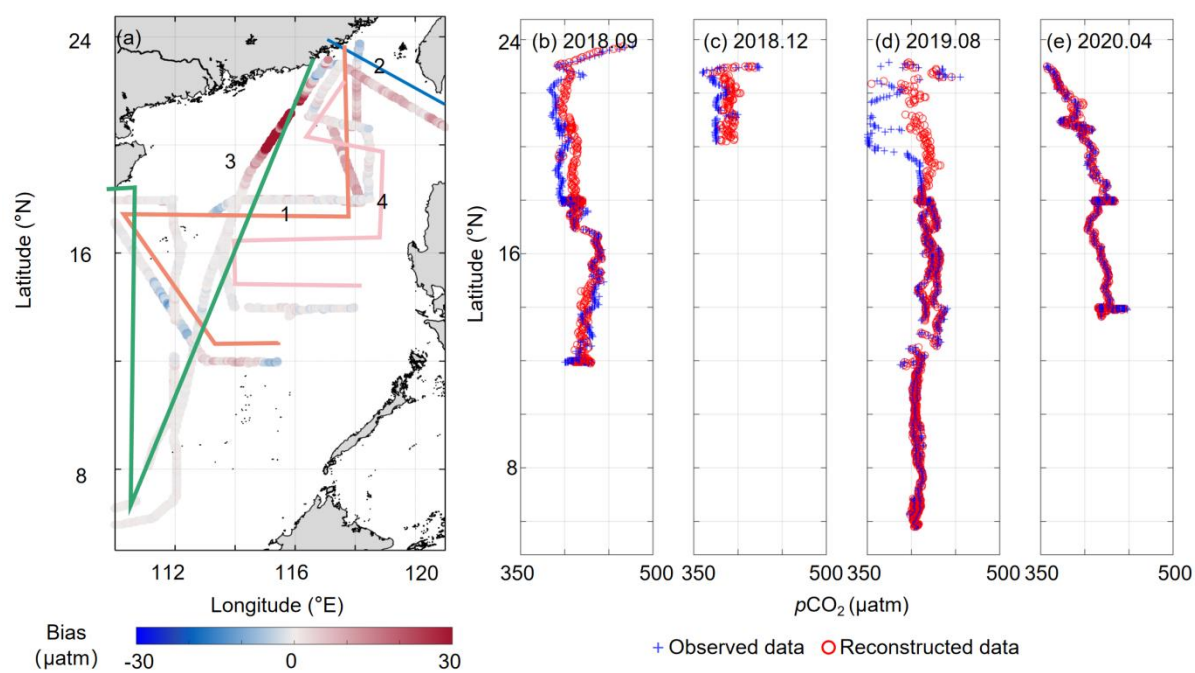
删除[Author]: no

删除[Author]: e

删除[Author]: reconstruction model



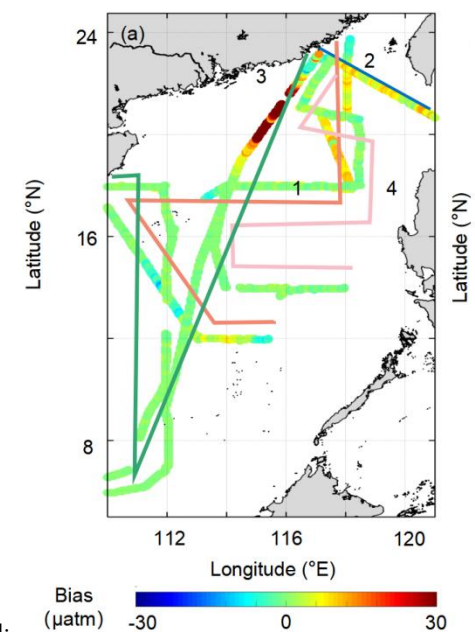
precipitation (~30mm/day; <https://psl.noaa.gov/data/gridded/data.ncep.reanalysis2.surface.html>) occurred around the Pearl River estuary region (24–20 °N), which led to intensification of the Pearl River plume. The plume has relatively low  $p\text{CO}_2$  values that eventually decreased the observed values along the coast. However, the monthly average runoff of the Pearl River during that month (August, 2019; <http://www.pearlwater.gov.cn/>; Pearl River Plume Index in Wang et al., 2022) was low, indicating that our retrieval algorithm is still highly reliable from the perspective of monthly averages. Thus, the inconsistencies between the reconstructed (monthly average) and the in situ datasets are mainly due to the differences in the time scales of the remote sensing and the in situ data. The reconstructed data in this study were determined on a monthly scale, while the temporal resolution of the in situ data was on the order of hours. It is clear that relatively pronounced short-term changes in  $p\text{CO}_2$ , such as the diurnal variability caused by short-term heavy precipitation, cannot be reflected in the reconstructed data.



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

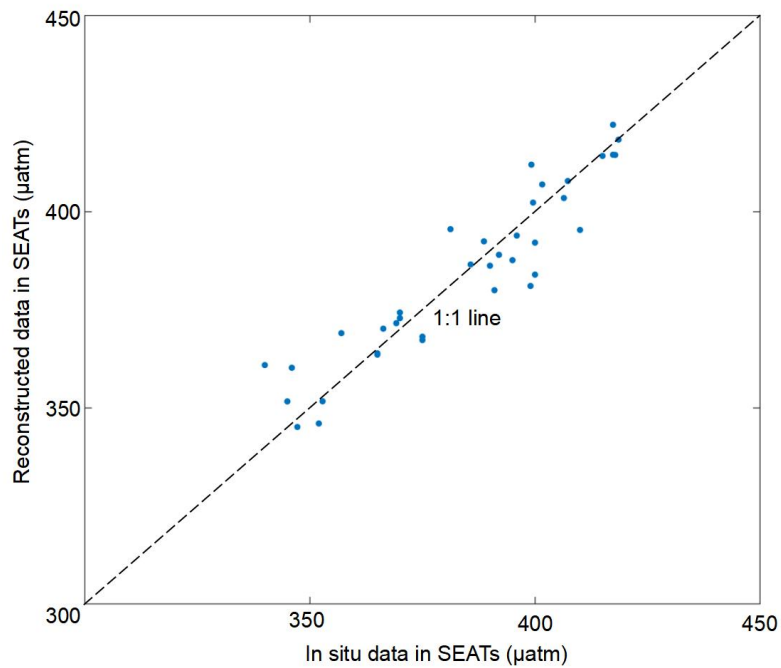
Dai et al. (2022) produced a time-series of in situ data from 2003 to 2019 at the SEATs station, which we used here to validate the accuracy of the long-term trends of our model data (results shown in Fig. 10). The long-term trend of reconstructed  $p\text{CO}_2$  data at the SEATs station is largely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

删除[Rick Smith]: ,  
 删除[Rick Smith]: such that the plume with  
 删除[Author]: in preparation  
 删除[Rick Smith]: is  
 删除[Author]: e  
 删除[Author]: reconstruction model  
 删除[Author]: a  
 删除[Rick Smith]: perspective  
 删除[Rick Smith]: inconsistency  
 删除[Author]: field-observed  
 删除[Rick Smith]: is  
 删除[Author]: field-observed  
 删除[Author]: field-observed  
 删除[Rick Smith]: variation



删除[Author]:  
 删除[Rick Smith]: ing  
 删除[Author]: observed  
 删除[Author]:  
 删除[Author]: observed data  
 删除[Author]:  
 删除[Rick Smith]: are  
 删除[Author]: observations  
 删除[Rick Smith]: of our





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

### 4.3 Uncertainties

As shown in Table 2, our reconstructed data have a high degree of accuracy, with an RMSE of  $\sim 10 \mu\text{atm}$  and MAE of  $\sim 6 \mu\text{atm}$ .

According to Equation 6, the bias of RS-derived  $p\text{CO}_2$  data used in the second term of Equation 6 is  $\sim 21 \mu\text{atm}$  (Table 2), the bias of SST is  $\sim 0.27^\circ\text{C}$  (Qin et al., 2014), the bias of SSS is  $\sim 0.33$  (Wang et al., 2022), and the bias of Chl-a is  $\sim 115\%$  (Zhang et al., 2006). We then estimated the  $p\text{CO}_2$  changes due to these features' variations by constraining these features based on our model, and computed  $\frac{\partial p\text{CO}_2}{\partial \text{Feature}}$ .

The overall uncertainty in the reconstructed dataset is greater in the coastal area ( $\sim 13 \mu\text{atm}$ ) than in the basin ( $\sim 10 \mu\text{atm}$ ) (Fig. 11 a), and this spatial pattern is mainly determined by the second term in Equation 6. The spatial distribution of the first term in Equation 6 (Fig. 11 b), calculated from a "max bias ratio," is consistent with that of  $p\text{CO}_2$  (Fig. 11 b). The second term in Equation 6 (Fig. 11 c) is calculated from the propagation of bias from each variable (Fig. 11 c). The Chl a bias (Fig. 11 f) shows it has the greatest effect on the reconstruction, among all the features (Fig. 11 f). Although the bias of the RS-derived  $p\text{CO}_2$  data is relatively large, the final influence it has on the results from the retrieval algorithm is negligible due to the use of the EOF method (Fig. 11 g).

删除[Rick Smith]: **the**

删除[Author]: **observations**

删除[Author]: **observed data**

删除[Author]: In previous studies, RMSE and MAE were

删除[Rick Smith]: reconstruction

删除[Rick Smith]: For the uncertainty

删除[Author]: calculations,

删除[Rick Smith]: a

删除[Author]: R1

删除[Rick Smith]:

删除[Author]: RS derived  $p\text{CO}_2$

删除[Author]: , However, this expression of uncertainty

删除[Rick Smith]: variability

删除[Author]: ing

删除[Author]: For example, for the  $\langle \text{math} \rangle$  part, we o

删除[Author]: results of uncertainty can be found in Fig.

删除[Author]: (Fig. 11 a)

删除[Rick Smith]: .

删除[Rick Smith]: A

删除[Rick Smith]: of

删除[Author]:  $p\text{CO}_2$

删除[Rick Smith]: of

删除[Rick Smith]: bias of

删除[Author]: between

删除[Rick Smith]: these

删除[Rick Smith]: .

删除[Author]: R

删除[Author]:

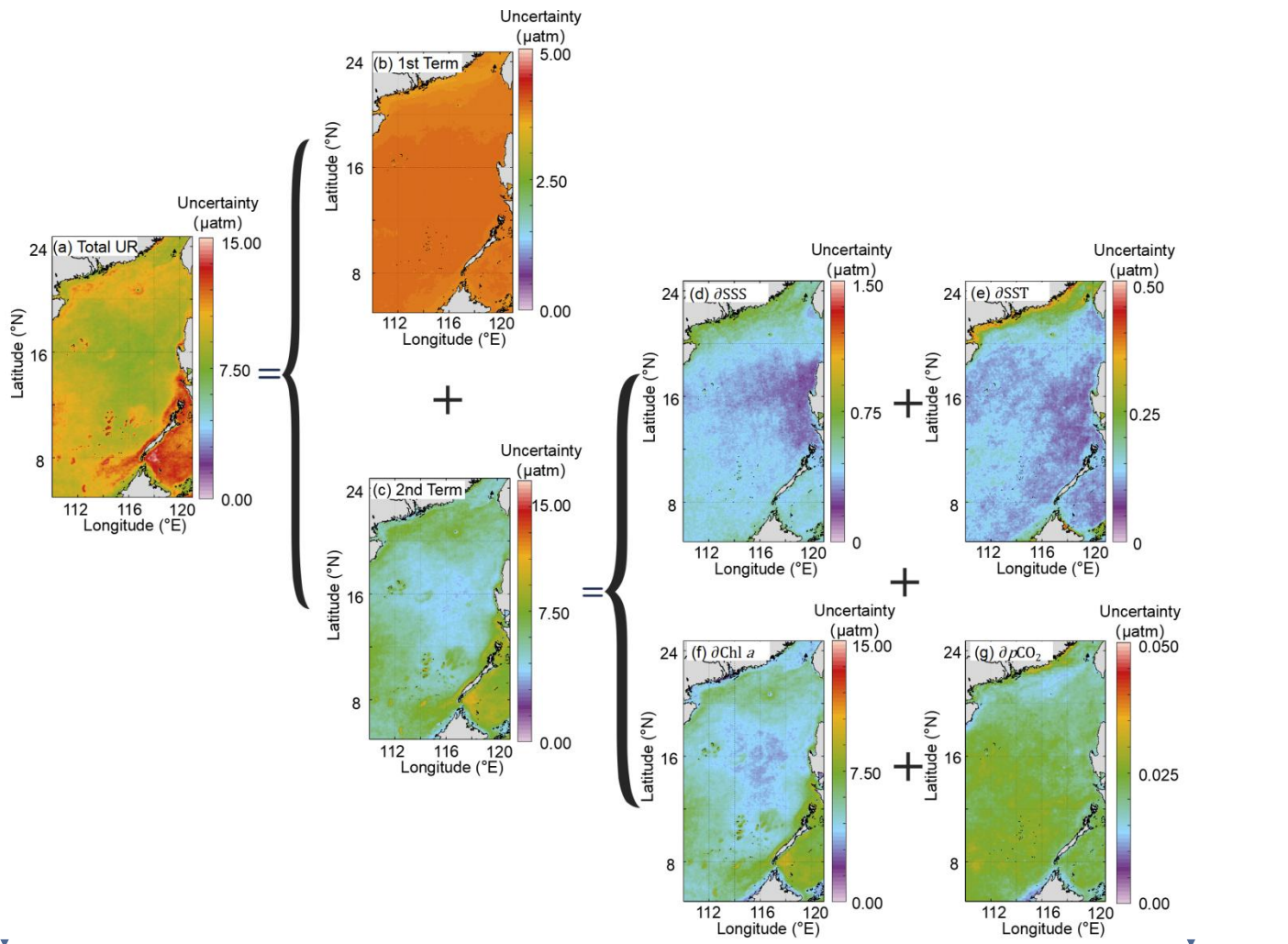
删除[Author]: S derived  $p\text{CO}_2$

删除[Rick Smith]: has

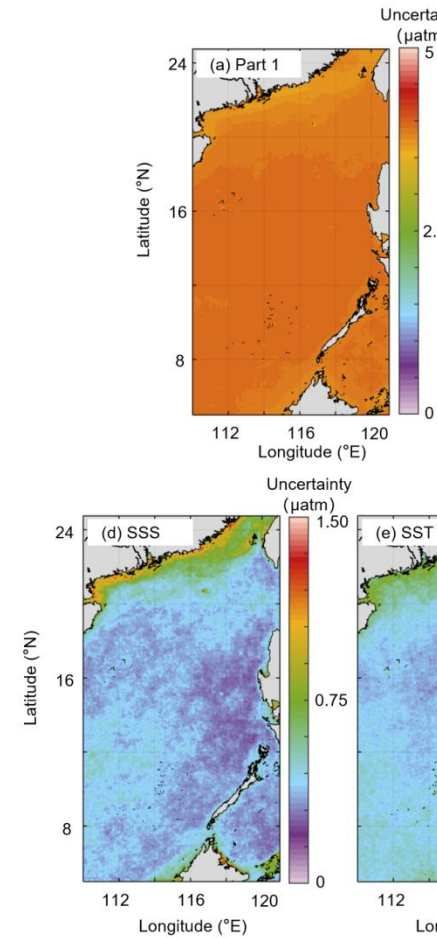
删除[Rick Smith]: bias

删除[Rick Smith]: of its bias

删除[Author]: retrieve algorithm reconstruction model



删除[Author]: These two parts were then added together



**Figure 11. Uncertainties of the reconstructed  $p\text{CO}_2$  fields (a, Total uncertainty in Equation 6; b, the first term of Equation 6; c, the second term of Equation 6; d,  $(\frac{\partial p\text{CO}_2}{\partial \text{SSS}})d\text{SSS}$  in the the second term of Equation 6; e,  $(\frac{\partial p\text{CO}_2}{\partial \text{SST}})d\text{SST}$  in the the second term of Equation 6; f,  $(\frac{\partial p\text{CO}_2}{\partial \text{Chl } a})d\text{Chl } a$  in the the second term of Equation 6; g,  $(\frac{\partial p\text{CO}_2}{\partial \text{RS derived } p\text{CO}_2})d\text{RS derived } p\text{CO}_2$  in the the second term of Equation 6.**

删除[Author]:

删除[Author]: stands for the

删除[Author]: stands for the

删除[Author]: stands for the

删除[Author]: stands for the

删除[Author]: <math>

删除[Author]: **Figure 11. Uncertainties of the**

删除[Rick Smith]: of the reconstructed  $p\text{CO}_2$  fields

删除[Rick Smith]: s

删除[Author]: in

设置格式[Rick Smith]: 非上标/ 下标

设置格式[Rick Smith]: 非上标/ 下标

删除[Rick Smith]: high-value

删除[Rick Smith]: center

删除[Rick Smith]: s

#### 4.4 Spatial and temporal $p\text{CO}_2$ features

The climatological monthly reconstructed  $p\text{CO}_2$  fields are shown in Figure 12. The highest values occur in May and June, and the lowest values occur in January. In winter,  $p\text{CO}_2$  first decreases in December and then increases after 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 increases from the basin to the northern coastal area, and the high  $p\text{CO}_2$  values in the central basin gradually expand outward starting in April. 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 consistently high values.



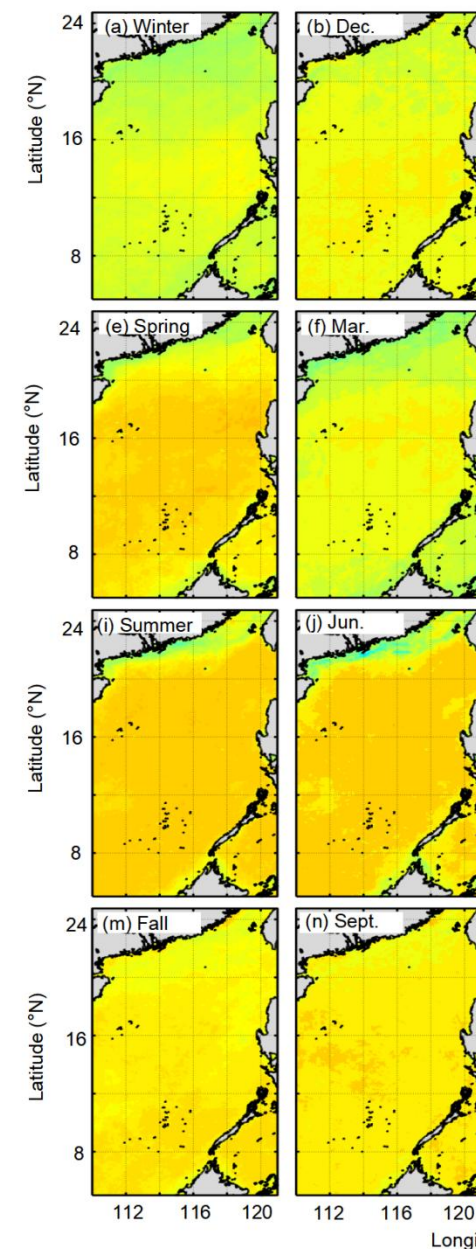
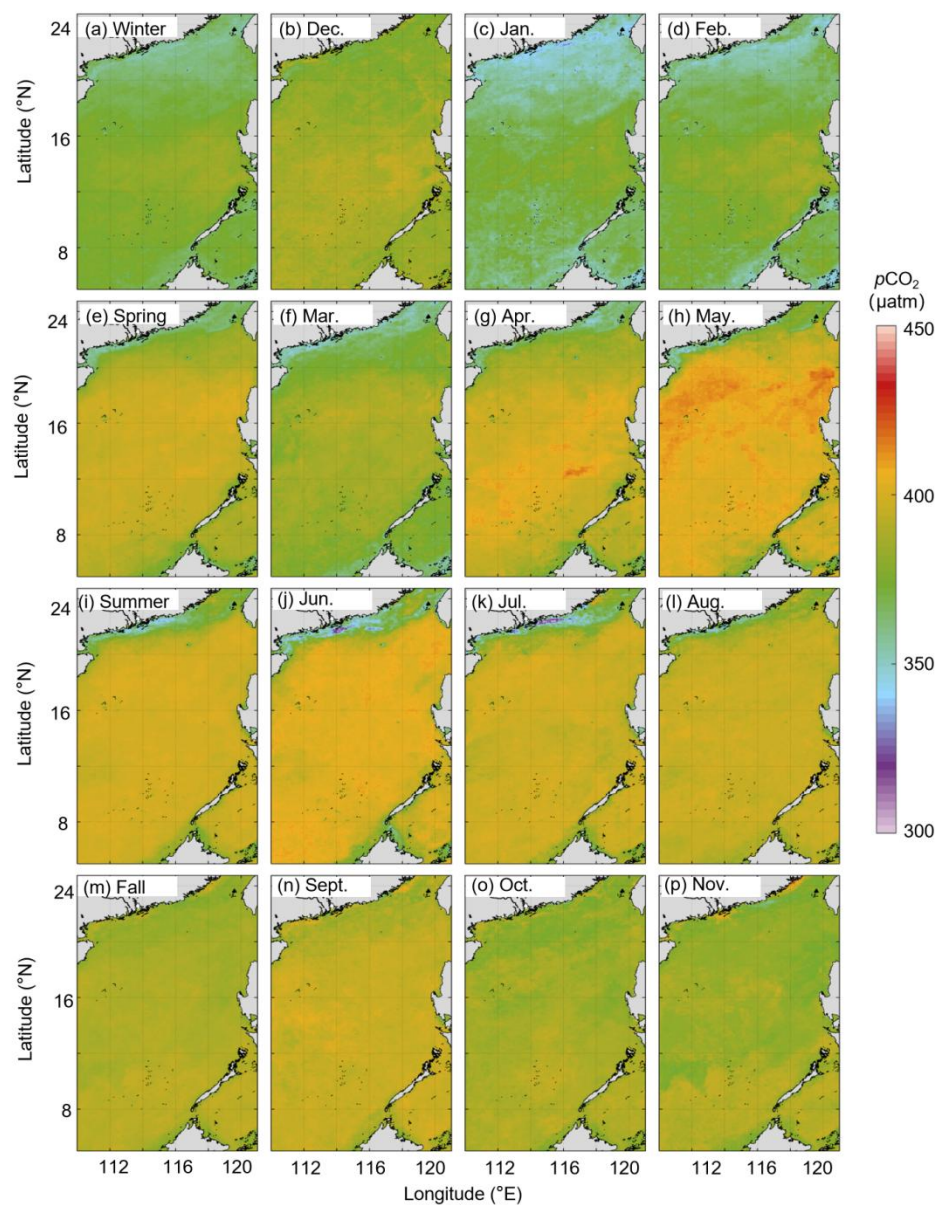


Figure 12. Long-term (2003–2020) seasonal and monthly average  $p\text{CO}_2$  field (unit:  $\mu\text{atm}$ ) (a. Winter; b. December; c. January; d. February; e. Spring; f. March; g. April; h. May; i. Summer; j. June; k. July; l. August; m. Fall; n. September; o. October; p. November).

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  $18^\circ\text{N}$  (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 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 controlled by the low SST, and the high  $p\text{CO}_2$  around Luzon Strait affected by winter upwelling. In summer, the reconstructed  $p\text{CO}_2$  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 effectively shows the influence of coastal upwelling in the northeastern shelf ( $\sim 23^\circ\text{N}$ ,  $117^\circ\text{E}$ ). Thus, our reconstructed  $p\text{CO}_2$  fields clearly reflect the spatial pattern of the *in situ*  $p\text{CO}_2$  (Fig. 3), which are generally consistent with previously reported patterns (Li et al., 2020; Zhai et al., 2013; Gan et al., 2010).

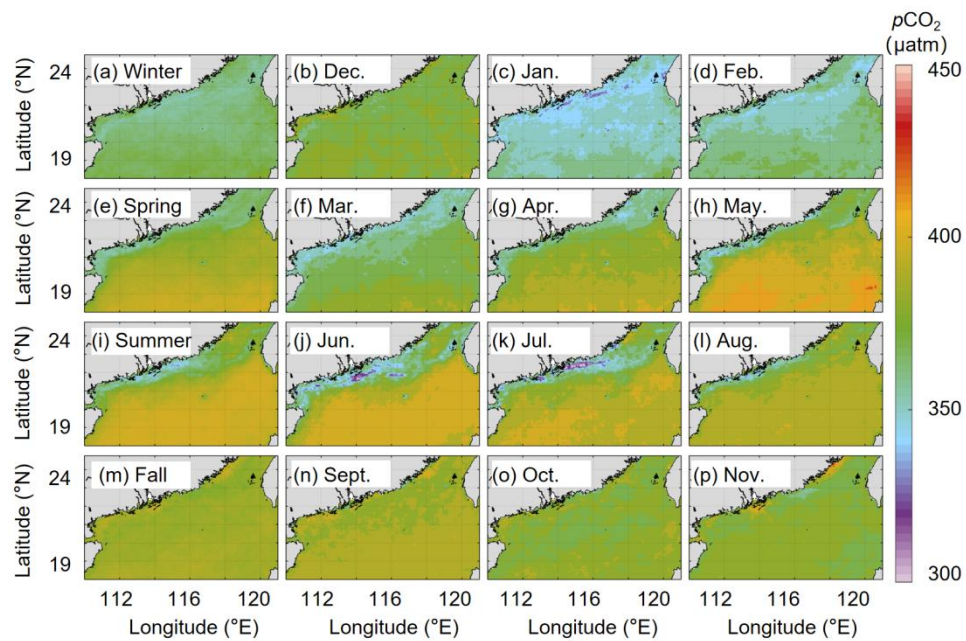
删除[Author]:

删除[Author]: w

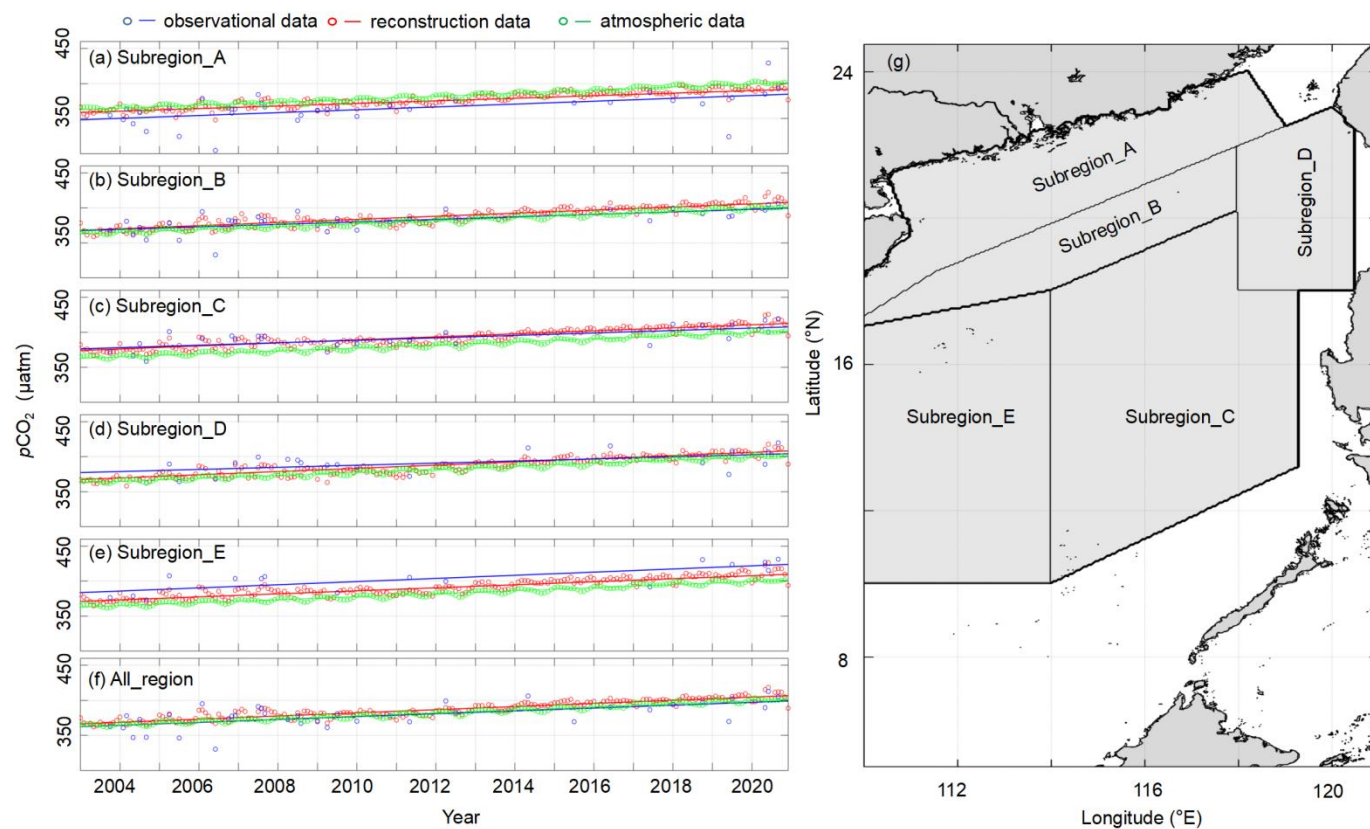
删除[Author]: .

删除[Author]: field observed

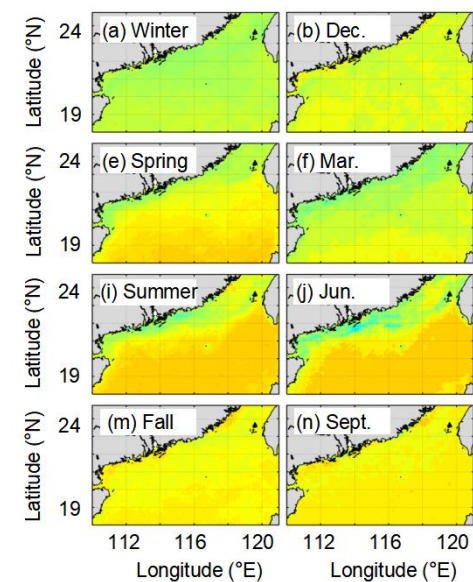




416  
417  
418 **Figure 13. Long-term (2003–2020) seasonal and monthly averaged  $p\text{CO}_2$  field in the region north of  $18^\circ\text{N}$  (unit:  $\mu\text{atm}$ ) (a.**  
419 **Winter; b. December; c. January; d. February; e. Spring; f. March; g. April; h. May; i. Summer; j. June; k. July; l. August;**  
420 **m. Fall; n. September; o. October; p. November).**



422  
423 **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)**



删除[Author]:

删除[Author]: w

删除[Author]: S

under study. The sub-regions are shown in (g). The lines indicate the deseasonalized long-term trend of the spatially averaged monthly  $p\text{CO}_2$  data for each sub-region with the slopes shown in Table 3. The deseasonalized method can be found in Landschützer et al. (2016).

**Table 4. Deseasonalized long-term trend of the spatially averaged monthly  $p\text{CO}_2$  data for each sub-region of the South China Sea. (unit:  $\mu\text{atm yr}^{-1}$ ).**

	All_region	Subregion_A	Subregion_B	Subregion_C	Subregion_D	Subregion_D
Reconstructed $p\text{CO}_2$	2.12±0.17	1.82±0.14	2.23±0.12	2.17±0.12	2.20±0.13	2.16±0.13
$\text{In situ } p\text{CO}_2$	2.10±0.79	1.80±0.86	1.73±0.84	1.81±0.85	1.41±1.16	2.13±1.10

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

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 the SCS is a weak source of atmospheric  $\text{CO}_2$ . This conclusion is consistent with previous studies (e.g., Li et al., 2020). Moreover, compared to the rate of atmospheric  $\text{CO}_2$  increase ( $\sim 2.2 \mu\text{atm yr}^{-1}$ ), for Subregion A, the  $p\text{CO}_2$  trend is much slower than that of 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 accumulation in this region is expected to increase in the future. For Subregion C and Subregion E, the spatially averaged monthly mean  $p\text{CO}_2$  is higher than the atmospheric  $p\text{CO}_2$ ; thus, these two regions will still provide a weak source of atmospheric  $\text{CO}_2$  in the future. Finally, whether Subregion B and Subregion D act as a source or sink of the atmospheric  $\text{CO}_2$  is influenced by seasonal changes and physical processes. Subregion B can be a zone of significant sink of atmospheric  $\text{CO}_2$  as demonstrated by its low sea surface  $p\text{CO}_2$  when the Pearl River plume spreads more widely in summer. In contrast, in winter when the Kuroshio intrusion is strong, both Subregions B and D have high sea surface  $p\text{CO}_2$ , indicating both subregions are sources of atmospheric  $\text{CO}_2$ .

## 5 Data availability

- 删除[Author]: ,
- 删除[Author]: .
- 删除[Author]: 3
- 删除[XHGuo]: *Data*
- 设置格式[XHGuo]: 字体: 倾斜
- 设置格式[XHGuo]: 下标
- 删除[Author]: Observation
- 删除[XHGuo]: Data
- 删除[Author]: region A
- 删除[Author]: region
- 删除[Author]: region
- 删除[Author]: region
- 删除[Author]: W
- 删除[Author]: region
- 删除[Author]: observational data
- 删除[Author]: observational data
- 设置格式[XHGuo]: 上标
- 删除[Author]: region
- 删除[Author]: egioms
- 设置格式[XHGuo]: 下标
- 删除[Author]: regions
- 删除[Author]: zone
- 删除[Author]: is spreading
- 删除[Author]: into a wider spatial coverage
- 删除[Author]: s
- 删除[Author]: When the Pearl River plume is relatively strong in summer, resulting in relatively low  $p\text{CO}_2$  in Sub\_region B, this sub\_region turns into a sink of atmospheric  $\text{CO}_2$ . When the Kuroshio invasion or water mixing is strong in winter, resulting in relatively high  $p\text{CO}_2$  in Sub\_region B and Sub\_region D, both two sub\_regional turn into a source of atmospheric  $\text{CO}_2$ .

452 The data (the reconstructed  $p\text{CO}_2$  data, the in situ  $p\text{CO}_2$  data before 2018 ( $0.5^\circ$  to  $0.5^\circ$ ), and the remote sensing derived  $\text{CO}_2$  data)  
453 for this paper are available under the link <https://doi.org/10.57760/sciencedb.02050> (Wang & Dai, 2022).

## 455 6 Conclusions

456 Based on the machine learning method, we reconstructed the sea surface  $p\text{CO}_2$  fields in the SCS with an  $0.05^\circ$  to  $0.05^\circ$  spatial  
457 resolution over the last two decades (2003-2020) by calculating the statistical relationship between the in situ  $p\text{CO}_2$  data and  
458 RS-derived data. The input data we used in machine learning include RS-derived data (sea surface salinity, sea surface  
459 temperature, chlorophyll), the spatial patterns of  $p\text{CO}_2$  calculated by EOF, atmospheric  $\text{CO}_2$ , and time labels (month). The  
460 machine learning method (CATBOOST) used in this study was facilitated by the EOF method, which provides spatial constraints  
461 for the data reconstruction. In addition to the typical machine learning performance metrics, we present a novel method for  
462 uncertainty calculation that incorporates the bias of both the reconstruction and the sensitivity of reconstructed models to its  
463 features. This method effectively shows the spatiotemporal patterns of bias, and makes up for the spatial representation of the  
464 typical performance metrics.

465 We validate our reconstruction with three independent testing datasets, and the results show that the bias between our  
466 reconstruction and in situ  $p\text{CO}_2$  data in the SCS is relatively small (about  $10 \mu\text{atm}$ ). Our reconstruction successfully captures the  
467 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  
468 further analyze the effect of meso-microscale processes (e.g., the Pearl River plume, and CCC) on sea surface  $p\text{CO}_2$  in the SCS.

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

473 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  
474 SCS and may be useful to predict changes of  $\text{CO}_2$  source or sink patterns in this system.

## 476 Author contribution

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

设置格式[XHGuo]: 字体: 倾斜

删除[Author]: Observational

设置格式[XHGuo]: 字体: 倾斜

删除[XHGuo]: \*

删除[Author]: <https://github.com/Elricriven/co2data>

删除[Author]: et al.,

删除[XHGuo]: ×

删除[admin]:

删除[Author]:

删除[admin]: high

删除[admin]: ( $0.05^\circ$  to  $0.05^\circ$ )

删除[Author]: underway observational

删除[admin]: remote sensing

删除[Author]:

删除[Author]: s

删除[admin]: remote sensing

删除[admin]: , because the latter can

删除[admin]: method

删除[Author]: observational

删除[admin]: shows

删除[Author]: s

删除[Author]: observational data

删除[Author]:  $p\text{CO}_2$



483 **Competing interests**

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

485

486 **Acknowledgements**

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

489

490

491 **References**

492 Bai, Y., Cai, W., He, X., Zhai, W., Pan, D., Dai, M., and Yu, P.: A mechanistic semi-analytical method for remotely sensing sea  
493 surface  $p\text{CO}_2$  in river-dominated coastal oceans: A case study from the East China Sea, *J. Geophys. Res.: Oceans*, 120,  
494 2331-2349, 2015.

495 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  
496 Ocean CO<sub>2</sub> Atlas (SOCAT), *Earth Syst. Sci. Data*, 8, 383-413, 2016.

497 Borges, A. V., Delille, B., and Frankignoulle, M.: Budgeting sinks and sources of CO<sub>2</sub> in the coastal ocean: Diversity of  
498 ecosystems counts, *Geophys. Res. Lett.*, 32, L14601, 2005.

499 Cao, Z. and Dai, M.: Shallow-depth CaCO<sub>3</sub> dissolution: Evidence from excess calcium in the South China Sea and its export to  
500 the Pacific Ocean, *Global Biogeochem. Cy.*, 25, GB2019, 2011.

501 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  
502 continental shelf system under the influence of both a river plume and coastal upwelling, *J. Geophys. Res.: Biogeo.*, 116,  
503 G02010, 2011.

504 Cao, Z., Yang, W., Zhao, Y., Guo, X., Yin, Z., Du, C., Zhao, H., and Dai, M.: Diagnosis of CO<sub>2</sub> dynamics and fluxes in global  
505 coastal oceans, *Natl. Sci. Rev.*, 7, 786-797, 2020.

506 Chen, C. and Borges, A. V.: Reconciling opposing views on carbon cycling in the coastal ocean: Continental shelves as sinks and  
507 near-shore ecosystems as sources of atmospheric CO<sub>2</sub>, *Deep-Sea Res. I*, 56, 578-590, 2009.

508 Chen, C., Lai, Z., Beardsley, R. C., Xu, Q., Lin, H., and Viet, N. T.: Current separation and upwelling over the southeast shelf  
509 of Vietnam in the South China Sea, *J. Geophys. Res.: Oceans*, 117, C03033, 2012.

510 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:  
511 Implications for understanding regional variability and air-sea gas exchange, *Geophys. Res. Lett.*, 34, 265-278, 2007.

512 Cheng, C., Xu, P. F., Cheng, H., Ding, Y., Zheng, J., Ge, T., and Xu, J.: Ensemble learning approach based on stacking for  
513 unmanned surface vehicle's dynamics. *Ocean Eng.*, 207, 107388, 2020.

514 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  
515 of atmospheric CO<sub>2</sub>?, *Geophys. Res. Lett.*, 40, 2154-2158, 2013.

516 Dai, M., Gan, J., Han, A., Kung, H., and Yin, Z.: “Physical Dynamics and Biogeochemistry of the Pearl River Plume” in  
517 *Biogeochemical Dynamics at Large River-Coastal Interfaces*. Eds. T. Bianchi, M. Allison and W. J. Cai (Cambridge University  
518 Press, Cambridge), 321-352, 2014.

519 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,  
520 G., and Wang, Z.: Carbon fluxes in the coastal ocean: Synthesis, boundary processes and future trends, *Annu. Rev. Earth Pl. Sc.*,  
521 50, 593-626, 2022.

522 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  
523 South China Sea: insights from an isopycnal mixing model, *Biogeosciences*, 10, 6419-6432, 2013.

524 Dong, L., Su, J. Wong, L. Cao, Z. and Chen, J.: Seasonal variation and dynamics of the Pearl River plume, *Cont. Shelf*  
525 *Res.*, 24, 1761-1777, 2004.

526 [Dye, A. W., Rastogi, B., Clemesha, R. E. S., Kim, J. B., Samelson, R. M., Still, C. J., & Williams, A. P.: Spatial patterns and trends](#)  
527 [of summertime low cloudiness for the Pacific Northwest, 1996–2017. \*Geophysical Research Letters\*, 47, e2020GL088121,](#)  
528 [2020.](#)

529 Fassbender, A. J., Rodgers, K. B., Palevsky, H. I., and Sabine, C. L.: Seasonal Asymmetry in the Evolution of Surface Ocean  
530 *pCO<sub>2</sub>* and pH Thermodynamic Drivers and the Influence on Sea - Air CO<sub>2</sub> Flux, *Global Biogeochem. Cy.*, 32, 1476-1497,  
531 2018.

532 Fay, A., Gregor, L., Landschützer, P., McKinley, G., Gruber, N., Gehlen, M., Iida, Y., Laruelle, G., Rödenbeck, C., Roobaert, A.,  
533 and Zeng, J.: SeaFlux: harmonization of air-sea CO<sub>2</sub> fluxes from surface *pCO<sub>2</sub>* data products using a standardized approach,  
534 *Earth Syst. Sci. Data*, 13, 4693-4710, 2021

535 Friedlingstein, P., Jones, M. W., O'Sullivan, M., Andrew, R. M., Hauck, J., Peters, G. P., Peters, W., Pongratz, J., Sitch, S., Le  
536 Quéré, C., Bakker, D. C. E., Canadell, J. G., Ciais, P., Jackson, R. B., Anthoni, P., Barbero, L., Bastos, A., Bastrikov, V., Becker,  
537 M., Bopp, L., Buitenhuis, E., Chandra, N., Chevallier, F., Chini, L. P., Currie, K. I., Feely, R. A., Gehlen, M., Gilfillan, D.,  
538 Gkritzalis, T., Goll, D. S., Gruber, N., Gutekunst, S., Harris, I., Haverd, V., Houghton, R. A., Hurtt, G., Ilyina, T., Jain, A. K.,  
539 Joetzjer, E., Kaplan, J. O., Kato, E., Klein Goldewijk, K., Korsbakken, J. I., Landschützer, P., Lauvset, S. K., Lefèvre, N.,  
540 Lenton, A., Lienert, S., Lombardozi, D., Marland, G., McGuire, P. C., Melton, J. R., Metzl, N., Munro, D. R., Nabel, J. E. M.  
541 S., Nakaoka, S.-I., Neill, C., Omar, A. M., Ono, T., Peregón, A., Pierrot, D., Poulter, B., Rehder, G., Resplandy, L., Robertson,  
542 E., Rödenbeck, C., Séférian, R., Schwinger, J., Smith, N., Tans, P. P., Tian, H., Tilbrook, B., Tubiello, F. N., van der Werf, G. R.,  
543 Wiltshire, A. J., and Zaehle, S.: Global Carbon Budget 2019, *Earth Syst. Sci. Data*, 11, 1783-1838, 2019.

544 Gan, J., Li, H., Curchitser, E. N., and Haidvogel, D. B.: Modeling South China sea circulation: Response to seasonal forcing

删除[Author]:

删除[Author]: pCO<sub>2</sub>

545 regimes, *J. Geophys. Res.: Oceans*, 111, C06034, 2006.

546 Gan, J., Li, L., Wang, D., and Guo, X.: Interaction of a river plume with coastal upwelling in the northeastern South China Sea,  
547 *Cont. Shelf Res.*, 29, 728-740, 2009.

548 Gan, J., Lu, Z., Dai, M., Cheung, A. Y. Y., Liu, H., and Harrison, P.: Biological response to intensified upwelling and to a river  
549 plume in the northeastern South China Sea: A modeling study, *J. Geophys. Res.: Oceans*, 115, C09001, 2010.

550 Guo, X. and Wong, G.: Carbonate chemistry in the Northern South China Sea shelf-sea in June 2010, *Deep Sea Res. II*, 117,  
551 119-130, 2015.

552 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  
553 transport from the East China Sea as a major nutrient source supporting winter primary production on the northeast South China  
554 Sea shelf, *Biogeosciences*, 10, 8159-8170, 2013.

555 Hu, J., Kawamura, H., Li, C., Hong, H., and Jiang, Y.: Review on current and seawater volume transport through the Taiwan Strait,  
556 *J. Oceanogr.*, 66, 591-610, 2010.

557 Jones, S. D., Quéré, C., and Rödenbeck, C.: Spatial decorrelation lengths of surface ocean  $f\text{CO}_2$  results in NetCDF format, *Global  
558 Biogeochem. Cy.*, 26, GB2042, 2014.

559 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  
560 Remote Sensing Based Neural Network Approach, *J. Geophys. Res.: Oceans*, 117, C08022, 2012. | 删除[Author]:  $p\text{CO}_2$

561 Landschützer, P., Gruber, N., and Bakker, D.: Decadal variations and trends of the global ocean carbon sink, *Global Biogeochem.  
562 Cy.*, 30, 1396-1417, 2016.

563 Landschützer, P., Gruber, N., and Bakker, D. C. E.: An updated observation-based global monthly gridded sea surface  $p\text{CO}_2$  and  
564 air-sea  $\text{CO}_2$  flux product from 1982 through 2015 and its monthly climatology, *Dataset*, 2017.

565 Laruelle, G., Lauerwald, R., Pfeil, B., and Regnier, P.: Regionalized global budget of the  $\text{CO}_2$  exchange at the air-water interface  
566 in continental shelf seas, *Global Biogeochem. Cy.*, 28, 1199-1214, 2015.

567 Landschützer, P., Bakker, D. C. E., Gruber, N., and Schuster, U.: Recent variability of the global ocean carbon sink, *Global  
568 Biogeochem. Cy.*, 28, 927-949, 2014.

569 Lefèvre, N., Watson, A., and Waston, A.: A comparison of multiple regression and neural network techniques for mapping in situ  
570  $p\text{CO}_2$  data, *Tellus B*, 57, 375-384, 2005.

571 Lefèvre, N., Watson, A. J., and Watson, A. R.: A comparison of multiple regression and neural network techniques for mapping in  
572 situ  $p\text{CO}_2$  data, *Tellus B: Chemical and Physical Meteorology, Dataset*, 2017.

573 [Levitus, S., Antonov, J. I., Boyer, T. P., Garcia, H. E., and Locarnini, R. A.: EOF analysis of upper ocean heat content,  
574 1956–2003, \*Geophys. Res. Lett.\*, 32, L18607, 2005.](#)

575 Li, Y., Xie, P., Tang, Z., Jiang, T., and Qi, P.: SVM-Based Sea-Surface Small Target Detection: A False-Alarm-Rate-Controllable



576 Approach, IEEE Geosci. Remote Sens., 16, 1225-1229, 2019.

577 Li, H., Wiesner, M. G., Chen, J., Lin, Z., Zhang, J., and Ran, L.: Long-term variation of mesopelagic biogenic flux in the central  
578 South China Sea: Impact of monsoonal seasonality and mesoscale eddy, Deep Sea Res. I, 126, 62-72, 2017.

579 Li, Q., Guo, X., Zhai, W., Xu, Y., Dai, M.: Partial pressure of CO<sub>2</sub> and air-sea CO<sub>2</sub> fluxes in the South China Sea: Synthesis of an  
580 18-year dataset, Prog. Oceanogr., 182, 102272, 2020.

581 Luo, X., Hao, W., Zhe, L., and Liang, Z.: Seasonal variability of air-sea CO<sub>2</sub> fluxes in the Yellow and East China Seas: A case  
582 study of continental shelf sea carbon cycle model, Cont. Shelf Res., 107, 69-78, 2015.

583 [McMonigal, K., & Larson, S. M.: ENSO explains the link between Indian Ocean dipole and Meridional Ocean heat  
584 transport. Geophysical Research Letters, 49, e2021GL095796, 2022.](#)

585 Mongwe, N. P., Chang, N., and Monteiro, P.: The seasonal cycle as a mode to diagnose biases in modelled CO<sub>2</sub> fluxes in the  
586 Southern Ocean, Ocean Model., 106, 90-103, 2016.

587 Park, J. H.: Effects of Kuroshio intrusions on nonlinear internal waves in the South China Sea during winter, J. Geophys. Res.:  
588 Oceans, 118, 7081-7094, 2013.

589 Qin, H., Chen, G., Wang, W., Wang, D., and Zeng, L.: Validation and application of MODIS-derived SST in the South China Sea,  
590 Int. J. Remote Sens., 35, 4315-4328, 2014.

591 Rödenbeck, C., Bakker, D. C. E., Gruber, N., Iida, Y., Jacobson, A. R., Jones, S., Landschützer, P., Metzl, N., Nakaoka, S., and  
592 Olsen, A.: Data-based estimates of the ocean carbon sink variability-first results of the Surface Ocean pCO<sub>2</sub> Mapping  
593 intercomparison (SOCOM), Biogeosciences, 12, 14-49, 2015.

594 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<sub>2</sub> flux at the  
595 SEATs time-series site, northern South China Sea, J. Geophys. Res.: Oceans, 115, C10021, 2010.

596 Tahata, M., Sawaki, Y., Ueno, Y., Nishizawa, M., Yoshida, N., Ebisuzaki, T., Komiya, T., and Maruyama, S.: Three-step  
597 modernization of the ocean: Modeling of carbon cycles and the revolution of ecological systems in the Ediacaran/Cambrian  
598 periods, Geosci. Front., 6, 121-136, 2015.

599 Telszewski, M., Chazottes, A., Schuster, U., Watson, A. J., and Wanninkhof, R.: Estimating the monthly pCO<sub>2</sub> distribution in the  
600 North Atlantic using a self-organizing neural network, Biogeosciences, 6, 1405-1421, 2009.

601 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  
602 basin-scale sea surface partial pressure of carbon dioxide from sparse in situ observations over the South China Sea, Earth Syst.  
603 Sci. Data, 13, 1403-1417, 2021.

604 Wanninkhof, R., Park, G. H., Takahashi, T., Sweeney, C., Feely, R., Nojiri, Y., Gruber, N., Doney, S. C., Mckinley, G. A., and  
605 Lenton, A.: Global ocean carbon uptake: magnitude, variability and trend, Biogeosciences, 10, 1983-2000, 2013

606 Wang, Z. and Dai, M.: [Datasets of reconstructed sea surface pCO<sub>2</sub> in the South China Sea, Science Data Bank](#),

删除[Author]: , Wang, G., Guo, X., Bai, Y., Xu, Y.,

删除[Author]: Spatial reconstruction of long-term (2003-2020)  
sea surface pCO<sub>2</sub> in the South China Sea using a machine  
learning based regression method aided by empirical  
orthogonal function analysis

删除[Author]: Github

- 607 <https://doi.org/10.57760/sciencedb.02050>.
- 608 [Wang, Z., Wang, G., Guo, X., Hu, J., and Dai, M. Reconstruction of High-Resolution Sea Surface Salinity over 2003–2020 in the](#)
- 609 [South China Sea Using the Machine Learning Algorithm LightGBM Model. \*Remote. Sens.\*, 14, 6147, 2022.](#)
- 610 <https://doi.org/10.3390/rs14236147>.
- 611 Xu, X., Zang, K., Zhao, H., Zheng, N., Huo, C., and Wang, J.: Monthly CO<sub>2</sub> at A4HDYD station in a productive shallow marginal  
612 sea (Yellow Sea) with a seasonal thermocline: Controlling processes, *J. Marine Syst.*, 159, 89-99, 2016.
- 613 Jo, Y., Dai, M., Zhai, W., Yan, X., and Shang, S.: On the variations of sea surface pCO<sub>2</sub> in the northern South China Sea: A remote  
614 sensing based neural network approach, *J. Geophys. Res.: Oceans*, 117, C08022, 2012.
- 615 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  
616 system under complex circulation schemes on a large continental shelf: The northern South China Sea, *Prog Oceanogr.*, 197,  
617 1026-1045, 2021.
- 618 [Yu, S., Song, Z., Bai, Y., and He, X.: Remote Sensing based Sea Surface partial pressure of CO<sub>2</sub> \(pCO<sub>2</sub>\) in China Seas](#)
- 619 [\(2003-2019\) \(2.0\). Zenodo, 2022. <https://doi.org/10.5281/zenodo.7372479>.](#)
- 620 Yu, Z., Shang, S., Zhai, W., and Dai, M.: Satellite-derived surface water pCO<sub>2</sub> and air-sea CO<sub>2</sub> fluxes in the northern South China  
621 Sea in summer, *Prog. Nat. Sci.*, 19, 775-779, 2009.
- 622 Zeng, J., Matsunaga, T., Saigusa, N., Shirai, T., Nakaoka, S. I., and Tan, Z. H.: Technical note: Evaluation of three machine  
623 learning models for surface ocean CO<sub>2</sub> mapping, *Ocean Sci.*, 13, 303-313, 2017.
- 624 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  
625 South China Sea in spring, summer and fall, *Mar. Chem.*, 96, 87-97, 2005.
- 626 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  
627 variations of sea-air CO<sub>2</sub> fluxes in the largest tropical marginal sea (South China Sea) based on multiple-year underway  
628 measurements, *Biogeosciences*, 10, 7775-7791, 2013.
- 629 Zhang, C., Hu, C., Shang, S., Müller-Karger, F., Yan, L., Dai, M., Huang, B., Ning, X., and Hong, H.: Bridging between SeaWiFS  
630 and MODIS for continuity of chlorophyll-a concentration assessments off Southeastern China, *Remote Sens. Environ.*, 102,  
631 250-263, 2006.
- 632 Zhan, Y., Zhang, H., Li, J., and Li, G.: Prediction Method for Ocean Wave Height Based on Stacking Ensemble Learning Model. *J.*  
633 *Mar. Sci. Eng.*, 10, 1150, 2022.
- 634 Zhu, Y., Shang, S., Zhai, W., and Dai, M.: Satellite-derived surface water pCO<sub>2</sub> and air-sea CO<sub>2</sub> fluxes in the northern South China  
635 Sea in summer, *Prog. Nat. Sci.*, 19, 775-779, 2009.

删除[Author]: <https://github.com/Elricriven/co2data>

删除[Author]: