



1     **MLAWind: A Monthly Sea Surface Wind Dataset Derived from an**  
2     **Interpretable Machine Learning Approach Integrating In-Situ**  
3     **Observations and Satellite Data**

4             Weihaio Guo<sup>1,2,3</sup>, Rongwang Zhang<sup>1,2,3</sup>, Xin Wang<sup>1,2,3</sup>, Dongxiao Wang<sup>4,5</sup>

6     **Affiliations**

7     <sup>1</sup>State Key Laboratory of Tropical Oceanography, South China Sea Institute of Oceanology, Chinese  
8     Academy of Sciences, Guangzhou, 510301, China.

9     <sup>2</sup>Global Ocean and Climate Research Center, South China Sea Institute of Oceanology, Chinese Academy  
10    of Sciences, Guangzhou, 510301, China.

11    <sup>3</sup>Guangdong Key Laboratory of Ocean Remote Sensing and Big Data, South China Sea Institute of  
12    Oceanology, Chinese Academy of Sciences, Guangzhou, 510301, China.

13    <sup>4</sup>School of Marine Sciences, Sun Yat-Sen University, Zhuhai, 519082, China.

14    <sup>5</sup>Southern Marine Science and Engineering Guangdong Laboratory (Zhuhai), Zhuhai, 519080, China.

15

16    *Correspondence to:* Rongwang Zhang (rwzhang@scsio.ac.cn); Xin Wang  
17    (wangxin@scsio.ac.cn)

18



19 **Abstract.** A gridded sea surface wind dataset with long temporal coverage is  
20 crucial for understanding atmospheric circulation changes and air-sea  
21 interactions at different time scales. This study employs an interpretable  
22 machine learning model based on random forest algorithm to generate a  $1^\circ \times 1^\circ$   
23 monthly sea surface wind dataset (MLAWind) from 1950 to 2023, covering the  
24 near-global ocean within  $60^\circ\text{S}$ – $60^\circ\text{N}$ . The data reconstruction model integrates  
25 the Cross-Calibrated Multi-Platform (CCMP) satellite data and the spatially  
26 sparse long-term International Comprehensive Ocean-Atmosphere Data Set  
27 (ICOADS), exhibiting robust interpretability and generalization capability.  
28 Evaluations demonstrate that the MLAWind dataset exhibits better agreement  
29 with remote sensing observations than existing reanalysis datasets during the  
30 training period (1993–2022), while maintaining robust performance during the  
31 independent testing period in 2023. Moreover, the performance of MLAWind  
32 since 1950 is assessed across multiple time scales. Its characteristics in  
33 climatology, annual cycle, and inter-annual variability are comparable to those  
34 of existing reanalysis datasets, even during the non-satellite period prior to  
35 1993. Uncertainties remain in the long-term trends of different datasets. The  
36 trend derived from MLAWind is corroborated by independent coral records  
37 during 1950–1982, which demonstrates its strong capability in reconstructing  
38 historical sea surface wind variations. The results indicate that MLAWind  
39 serves as a reliable data resource for global climate change research. The  
40 reconstructed MLAWind dataset is publicly accessible at  
41 <https://doi.org/10.5281/zenodo.17354864> (Guo et al., 2025b).

42



## 43 **1. Introduction**

44 Sea surface wind is a critical factor in air-sea interactions, exerting  
45 significant impacts on climate change, marine ecosystem and human society  
46 (Tokinaga et al., 2012; Wang et al., 2018; Zhou et al., 2022). It governs  
47 variations of the hydrological cycle by modulating atmosphere-ocean heat and  
48 moisture exchanges (Held and Soden, 2006; Findell et al., 2019). Sea surface  
49 wind drives the large-scale redistribution of water masses and energy. The  
50 global ocean current systems, including the Kuroshio, Gulf Stream, South  
51 China Sea Throughflow, Indonesian Throughflow, and Antarctic Circumpolar  
52 Current, are primarily regulated by wind-driven Ekman transport and  
53 momentum transfer processes (Deser et al., 1999; Wang et al., 2006, 2023;  
54 Zhang et al., 2023; Li et al., 2025). Sea surface wind can deliver energy to the  
55 deep ocean, inducing near-inertial waves and near-bottom currents (Zhang et  
56 al., 2024). Previous studies have demonstrated that equatorial sea surface wind  
57 anomalies are associated with the inter-annual variability of the tropical sea  
58 surface temperature (SST) anomalies, which promotes the development of El  
59 Niño-Southern Oscillation (ENSO) events (Kuo, et al., 2009; Clarke, 2014;  
60 Wang, 2019). Sea surface wind is thus utilized as a precursor to improve the  
61 predictability of different types of ENSO (Ren et al., 2019; Tseng et al., 2022).

62 The existing observational sea surface wind datasets exhibit notable  
63 limitations despite remarkable advances in observation platforms and sensor  
64 technologies. Satellite data are characterized by exceptional spatiotemporal  
65 continuity, offering a robust data foundation for retrieving climate system  
66 characteristics at the global scale, including those associated with extreme  
67 meteorological events (Vinoth and Young, 2011; Yang et al., 2015; Young and  
68 Ribal, 2019). However, the relatively short temporal coverage of satellite  
69 observations constrains their capacity to analyze climate variability from inter-  
70 decadal time scale to long-term trends. The International Comprehensive  
71 Ocean-Atmosphere Data Set (ICOADS) is widely acknowledged as the most



72 extensive and comprehensive global surface ocean dataset. It provides in-situ  
73 observations of critical atmospheric and oceanic variables dating back to 1662  
74 (Freeman et al., 2017). In contrast to the extensive spatial coverage of satellite  
75 data, ICOADS exhibit an irregular distribution and suffer from significant  
76 sampling errors. Due to these deficiencies in existing observational datasets,  
77 the Sixth Assessment Report by the Intergovernmental Panel on Climate  
78 Change (IPCC AR6) assigns a 'low to medium' confidence level for the  
79 historical sea surface wind trends assessment (IPCC, 2021).

80 A long-term gridded global wind dataset can be constructed by integrating  
81 in-situ and satellite observations. Numerous reanalysis products have  
82 assimilated the ICOADS data and multi-platform satellite observations through  
83 data assimilation techniques, significantly extending the temporal-spatial  
84 coverage of the sea surface wind field (e.g. Kalnay et al., 1996; Kobayashi et  
85 al., 2015; Hersbach et al., 2019). The inherent limitations of reanalysis datasets  
86 primarily stem from two aspects: (1) the assimilation methodologies and  
87 models introduce significant biases into the reanalysis products; (2) variations  
88 in the types and volumes of in-situ data result in sampling biases. Substantial  
89 uncertainties persist in reanalysis datasets with different assimilation  
90 methodologies and data sources (Zhang et al., 2023).

91 In addition to data assimilation techniques, direct reconstruction methods  
92 based on ICOADS data are employed in some studies (Berry and Kent, 2011;  
93 Tokinaga et al., 2011, 2012). For example, the Wave- and Anemometer-Based  
94 Sea Surface Wind (WASWind) achieves high accuracy by integrating  
95 anemometer-measured winds and wave height-estimated winds from the  
96 ICOADS data (Tokinaga and Xie, 2011). It exhibits a relatively coarse spatial  
97 resolution of  $4^{\circ} \times 4^{\circ}$  and limited temporal coverage from 1950 to 2008. The  
98 National Oceanography Centre Southampton Flux Dataset v2.0 (NOCS2) is  
99 another reconstructed dataset derived from the ICOADS data. It features an  
100 enhanced spatial resolution of  $1^{\circ} \times 1^{\circ}$  but covers a shorter time period (1973–



101 2014) for sea surface wind speed (Berry and Kent, 2011).

102 Machine learning algorithms demonstrate superior capability in capturing  
103 nonlinear relationships between different variables compared to traditional  
104 techniques (Reichstein et al., 2019; Jiang et al., 2024; Wang and Li, 2024).  
105 These algorithms are applied in the reconstruction of various datasets such as  
106 sea surface temperature (Huang et al., 2025), air temperature (Yoo et al., 2018;  
107 He et al., 2022), ocean salinity (Tian et al., 2022), ocean heat content (Su et al.,  
108 2020; Bagnell and DeVries, 2021), carbon and water fluxes (Leng et al., 2024),  
109 and terrestrial water storage (Yin et al., 2023). Traditional machine learning  
110 models suffer from interpretability challenges, as their input-output mappings  
111 are treated as "black-box" systems lacking physically meaningful explanations.  
112 Recently, resolving machine learning interpretability issues has emerged as a  
113 critical research priority (Tian et al., 2022; Qin et al., 2024). SHapley Additive  
114 exPlanations (SHAP) is a reliable interpretable module based on cooperative  
115 game theory, which demonstrates excellent compatibility across diverse  
116 machine learning models. By generating both global and local interpretability,  
117 SHAP enables a more comprehensive attribution analysis than conventional  
118 feature importance metrics. The continuous advancement of interpretable  
119 machine learning techniques offers a robust methodology for historical sea  
120 surface wind reconstruction.

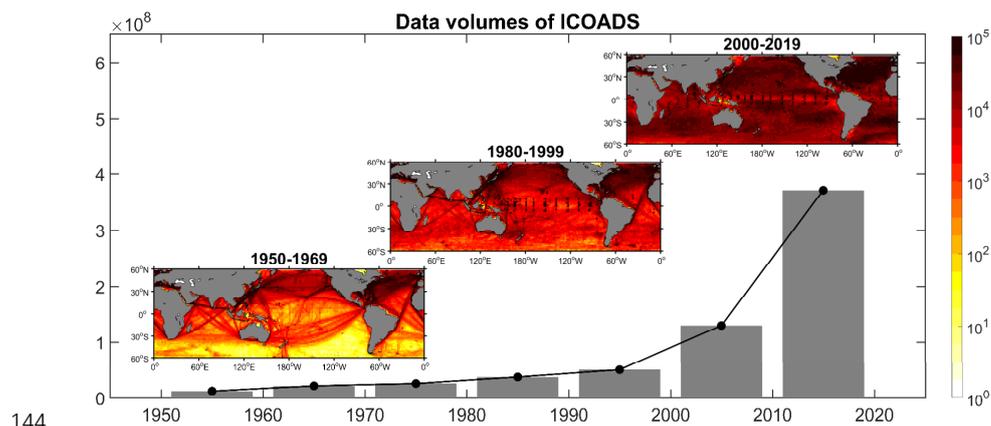
121 This study employs an interpretable machine learning algorithm to  
122 integrate ICOADS and satellite data, constructing a  $1^\circ \times 1^\circ$  monthly sea surface  
123 wind dataset from 1950–2023 with near-global ocean coverage within  $60^\circ\text{S}$ –  
124  $60^\circ\text{N}$ . The rest of the paper is organized as follows: the datasets, methods and  
125 data reconstruction model used in this study are described in Sect. 2; the  
126 evaluation processes of the reconstructed dataset are shown in Sect. 3; the  
127 performances of the reconstructed dataset at different time scales are examined  
128 in Sect. 4; the data availability statement is provided in Sect. 5; the study is  
129 summarized in Sect. 6.



## 130 2. Data and methods

### 131 2.1 Input data for model establishment

132 The reconstructed dataset is derived from in-situ observations of the  
133 ICOADS dataset. The ICOADS dataset provides an extensive collection of  
134 observations from diverse observing systems, containing multiple essential  
135 atmospheric and oceanic variables (Freeman et al., 2017). Sea surface wind speed  
136 (WS), wind direction (WD), SST, air temperature (Ta), and sea level pressure  
137 (SLP) from 1950 to 2023 are utilized in this study. Note that the ICOADS  
138 observations exhibit pronounced spatiotemporal heterogeneity. During the  
139 period 1950–1969, data coverage is relatively limited in the Southern  
140 Hemisphere, contrasting with denser distributions in the North Pacific and  
141 North Atlantic Oceans. From 1980 to 1999, the data volume gradually increases,  
142 followed by significant growth after 2000, which achieves near-global coverage  
143 (Fig. 1).



144  
145 **Figure 1.** The volumes of ICOADS dataset since 1950. The bar illustrates the decadal cumulative  
146 volume of ICOADS dataset. The spatial distributions are presented in different periods.

147 Following previous studies (Berry and Kent, 2011; Tokinaga and Xie,  
148 2011), comprehensive preprocessing is performed on the ICOADS dataset. All  
149 variables in the ICOADS dataset are subjected to rigorous quality control  
150 through an outlier detection procedure based on a climatological threshold of



151 3.5 standard deviations. The selection of the threshold aligns with the  
152 standardized processing implemented for both the ICOADS Monthly Summary  
153 Groups (MSG) product and the WASWind dataset (Tokinaga and Xie, 2011).  
154 Additional height adjustments are required for ICOADS sea surface winds, as  
155 the altitudes of ships and measurement instruments change over time. The bulk  
156 formulae and related parameters established by Smith (1980) are utilized to  
157 unify the unadjusted ICOADS sea surface winds to a standard reference height  
158 of 10 m. All data are initially processed into monthly means to enable point-to-  
159 point learning during model establishment process.

160 Satellite sea surface wind from the cross-calibrated multiplatform (CCMP)  
161 dataset is another important source of the reconstructed dataset. The CCMP  
162 dataset covers the period from 1993 to 2023 at a horizontal resolution of  
163  $0.25^{\circ} \times 0.25^{\circ}$  (Mears et al., 2019). It is spatially interpolated onto a  $1^{\circ} \times 1^{\circ}$  grid  
164 to facilitate point-to-point learning with ICOADS dataset. It integrates multiple  
165 advanced satellites, including Advanced Scatterometer-A and B (ASCAT-A/B),  
166 Special Sensor Microwave/Imager (SSM/I), Special Sensor Microwave  
167 Imager/Sounder (SSMIS), Tropical Rainfall Measuring Mission's (TRMM)  
168 Microwave Imager (TMI), Global Precipitation Measurement (GPM)  
169 Microwave Imager (GMI), Advanced Microwave Scanning Radiometer-EOS  
170 (ASMR-E), Advanced Microwave Scanning Radiometer 2 (AMSR2), Quick  
171 Scatterometer (QuikScat), and Wind Satellite (WindSat).

## 172 **2.2 Data for evaluation**

173 Five reanalysis products, including European Centre for Medium-Range  
174 Weather Forecasts (ERA5, Hersbach et al., 2019), the Japanese 55-year  
175 Reanalysis (JRA-55, Kobayashi et al., 2015), the National Oceanic and  
176 Atmospheric Administration 20th-Century Reanalysis (NOAA-20C, Compo et  
177 al., 2011), the National Centers for Environmental Prediction reanalysis 1  
178 (NCEP1, Kalnay et al., 1996) and its upgraded version (NCEP2, Kanamitsu et  
179 al., 2002), are compared with the reconstructed dataset. All reanalysis products



180 are spatially interpolated onto a  $1^\circ \times 1^\circ$  grid. Comprehensive descriptions of  
181 these datasets are presented in Table 1.

182 Three coral Mn/Ca records provided by *Porites* spp. corals are utilized to  
183 evaluate the long-term trends of sea surface wind. The coral record derived  
184 from Tarawa Atoll ( $1.33^\circ\text{N}$ ,  $172.97^\circ\text{E}$ ) is resampled at approximately quarterly  
185 (4 per year) increments, spanning from 1894 to 1982 (Thompson et al., 2015).  
186 The coral records from Butaritari Atoll ( $3.07^\circ\text{N}$ ,  $172.75^\circ\text{E}$ ) and Kiritimati  
187 Island ( $1.93^\circ\text{N}$ ,  $157.49^\circ\text{W}$ ) exhibit bimonthly temporal resolution, covering the  
188 periods of 1989–2010 and 1994–2011, respectively (Sayani et al., 2021). These  
189 islands are located far from continental Mn sources and possess a west-facing  
190 lagoon that is sheltered from the easterly trade winds. Their Mn/Ca records  
191 demonstrate a strong correlation with westerly wind events in the equatorial  
192 Pacific Ocean, serving as a valuable and independent resource for evaluating  
193 the zonal wind anomalies (Thompson et al., 2015; Sayani et al., 2021).

194 **Table 1.** A list of the datasets used in this study.

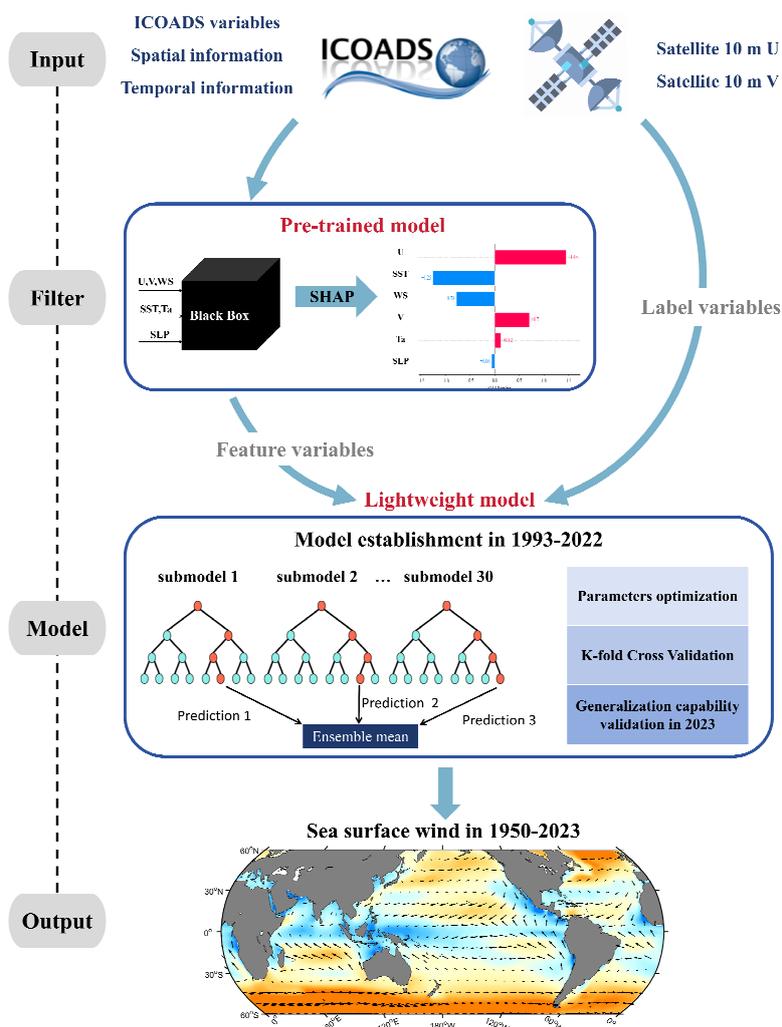
	Product	Variable	Data period	Resolution	Reference
In-situ observations	ICOADS	WS, WD, SST, Ta, SLP	1950–2023	-	Freeman et al. (2017)
Satellite data	CCMP	U, V, WS	1993–2023	$0.25^\circ \times 0.25^\circ$	Mears et al. (2019)
	ERA5	U, V, WS	1950–2022	$0.25^\circ \times 0.25^\circ$	Hersbach et al. (2019)
	NCEP1	U, V, WS	1950–2022	$1.875^\circ \times 1.90^\circ$	Kalnay et al. (1996)
	NCEP2	U, V, WS	1979–2022	$1.875^\circ \times 1.90^\circ$	Kanamitsu et al. (2002)
	JRA-55	U, V, WS	1958–2022	$1.25^\circ \times 1.25^\circ$	Kobayashi et al. (2015)
Reanalysis datasets	NOAA-20C	U, V, WS	1950–2014	$2^\circ \times 2^\circ$	Compo et al. (2011)
	Tarawa	Mn/Ca	1950–1982	-	Thompson et al. (2015)
	Butaritari	Mn/Ca	1989–2010	-	Sayani et al. (2021)
Corals	Kiritimati	Mn/Ca	1994–2011	-	Sayani et al. (2021)

### 195 2.3 Data reconstruction model

196 Figure 2 illustrates the schematic diagram of the construction process of  
197 the machine learning-assisted sea surface wind dataset (MLAWind). The  
198 random forest model is employed as the core algorithm. It effectively captures



199 nonlinear relationships between feature and label variables through multiple  
 200 decision trees (Breiman, 2001) and is extensively employed for temporal-  
 201 spatial data processing (Yoo et al., 2018; Watt-Meyer et al., 2021; He et al.,  
 202 2022). Notably, the inherent randomness in random forest model significantly  
 203 improves its generalization capability, which is critical for robust  
 204 reconstruction of long-term historical sea surface wind.



205  
 206 **Figure 2.** Schematic diagram of the construction process of machine learning-assisted sea surface  
 207 wind dataset (MLAWind).



208           The reconstruction framework comprises four essential steps. The first is  
209 tokenization and preprocessing of the input data. Multiple ICOADS variables  
210 (WS, WD, SST, Ta, SLP) are used as feature inputs for the data reconstruction  
211 model, while the CCMP 10-m wind field serves as the label variable. All feature  
212 values for the data reconstruction model are normalized to enhance the  
213 efficiency of parameter optimization (Su et al., 2018; Watt-Meyer et al., 2021).

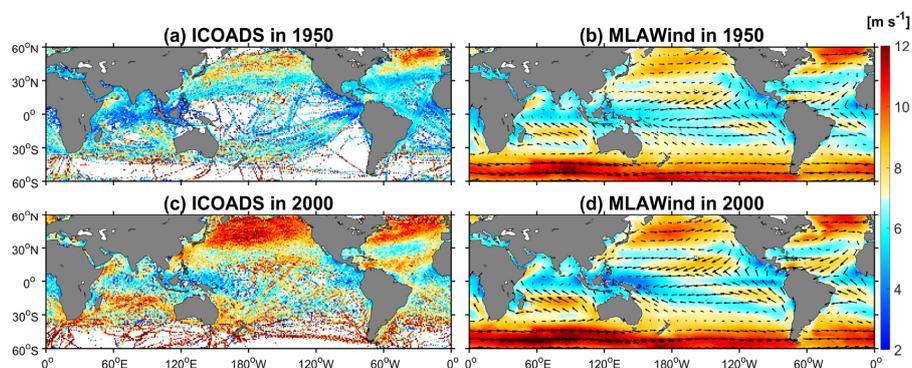
214           The second step is establishing a pre-trained model using the SHAP  
215 algorithm. SHAP is an interpretable artificial intelligence approach that  
216 quantifies the contribution of inputs and works as a filter to eliminate irrelevant  
217 features, thereby enhancing model interpretability (Lundberg and Lee, 2017).  
218 Based on the SHAP analysis, the ICOADS 10-m zonal wind, meridional wind,  
219 SST, WS, longitude, latitude, and month are identified as the seven key feature  
220 variables for sea surface wind reconstruction. They are utilized to establish a  
221 lightweight model, while other variables with little contribution are excluded.

222           Thirdly, the lightweight model is trained from 1993–2022, facilitating  
223 point-to-point learning between ICOADS variables and CCMP sea surface  
224 winds. Parameter tuning, K-fold Cross Validation, and generalization  
225 capability validation are three key components for lightweight models: (1) a  
226 grid search approach is employed to identify the optimal model parameters.  
227 Consequently, the number of regression trees, the minimum number of samples  
228 required to split and samples required to be at a leaf node is set to 200, 3, 3,  
229 respectively; (2) leave-one-out Cross Validation (LOOCV), a specialized form  
230 of K-fold Cross Validation is implemented during the model establishment  
231 process (Mao et al., 2014). The 30-year samples from 1993–2022 are segmented  
232 into 30 submodels, with each submodel consisting of a 29-year training set and  
233 a 1-year validation set. LOOCV exhibits minimal discrepancies among these  
234 model, indicating the lightweight model's robust performance; (3) an  
235 independent verification is conducted in 2023 to evaluate model's  
236 generalization capability. The results will be presented in Sect. 3.



237 The final step is to output the restructured dataset. Following the  
238 aforementioned training process conducted from 1993 to 2022, the model  
239 derives a nonlinear relationship between ICOADS variables and CCMP sea  
240 surface wind. The established relationship enables the reconstruction of  
241 historical sea surface winds during the period without satellite observations  
242 (1950–1992) by utilizing ICOADS data since 1950. The resultant monthly  
243 MLAWind dataset spanning 1950 to 2023, covering near-global oceans  
244 between 60°S and 60°N at a horizontal resolution of 1°×1°.

245 Figure 3 exhibits the reconstruction efficacy in 1950 and 2000. The two  
246 years are chosen as examples to represent periods with sparse and abundant  
247 original data coverage, respectively (Fig. 1). The original ICOADS from 1950  
248 exhibits significant spatial data gaps, particularly in the tropical Pacific and  
249 Southern Ocean. MLAWind shows continuous spatial coverage. Despite limited  
250 observations in the Southern Ocean, it demonstrates robust reconstruction of  
251 the prevailing westerly wind belt (Garreaud et al., 2013). In 2000, ICOADS  
252 exhibits a higher spatial density, while still containing considerable noise.  
253 MLAWind demonstrates smoother characteristics and better capability in  
254 capturing sea surface wind patterns.



255 **Figure 3.** Reconstruction results presentation. Sea surface wind of the original ICOADS and  
256 MLAWind dataset are compared in 1950 and 2000, respectively. These two years are chosen as  
257 examples to represent periods with sparse and abundant original data coverage, respectively.  
258

259



## 260 2.4 Evaluation metrics

261 Five metrics are employed in the evaluation process, including mean error  
262 (Bias), root mean square error (RMSE), Pearson correlation coefficient  
263 (CORR), skewness and the coefficient of determination (R-squared). Bias,  
264 RMSE and CORR quantify discrepancies between observations and products.  
265 Skewness characterizes the asymmetry of a variable's probability distribution.  
266 Positive values denote right-skewed distributions with longer tails toward large  
267 values and a left-shifted peak. R-square serves as a key metric for assessing  
268 model performance. Its values span from 0 to 1, with those closer to 1 indicating  
269 better model performance. These metrics can be described as follows:

$$270 \text{ Bias} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i) \quad (1)$$

$$271 \text{ RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (2)$$

$$272 \text{ CORR} = \frac{\sum_{i=1}^n (y_i - \bar{y})(\hat{y}_i - \bar{\hat{y}})}{\sqrt{\sum_{i=1}^n (y_i - \bar{y})^2} \sqrt{\sum_{i=1}^n (\hat{y}_i - \bar{\hat{y}})^2}} \quad (3)$$

$$273 \text{ Skewness} = \frac{\sqrt{n(n-1)}}{n-2} \frac{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^3}{\left(\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2\right)^{3/2}} \quad (4)$$

$$274 R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad (5)$$

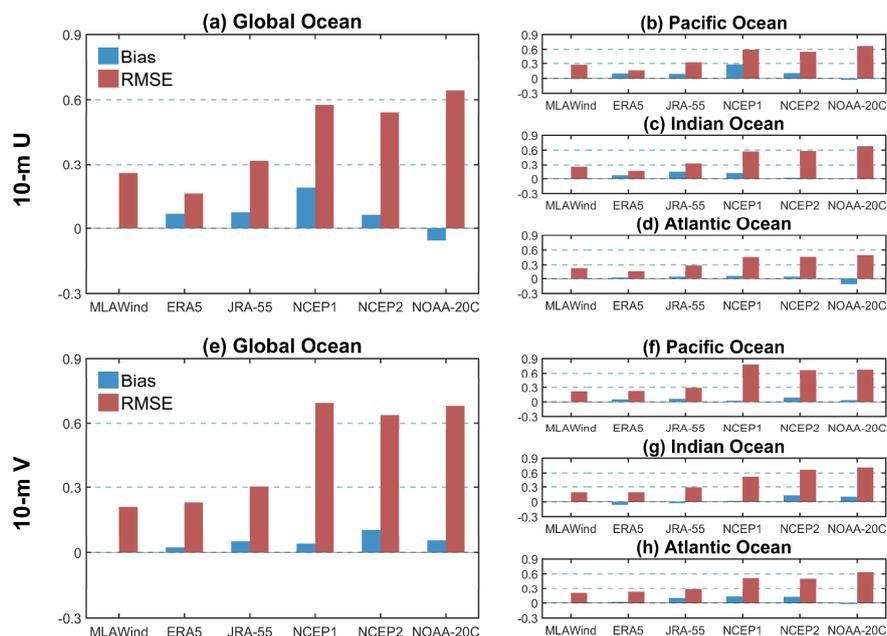
275 where  $n$  represents the number of samples,  $y_i$  represents the values obtained  
276 from different products,  $\hat{y}_i$  represents observational values.  $\bar{y}$  and  $\bar{\hat{y}}$  represent the  
277 mean values of  $y_i$  and  $\hat{y}_i$ , respectively.

## 278 3. Evaluation of the MLAWind dataset

279 Based on satellite data, we systematically evaluate the performance of the  
280 MLAWind dataset on the training set (1993–2022) and testing set (2023),  
281 respectively (Figs. 4-6). During the training period, the MLAWind dataset  
282 shows the smallest mean biases for both 10-m zonal and meridional winds (Fig.  
283 4a, e). It achieves the lowest RMSE ( $0.21 \text{ m s}^{-1}$ ) for meridional wind, while its  
284 RMSE ( $0.26 \text{ m s}^{-1}$ ) ranks second only to the ERA5 dataset ( $0.16 \text{ m s}^{-1}$ ) for zonal  
285 wind. The mean bias and RMSE are quantified across the Pacific, Indian, and

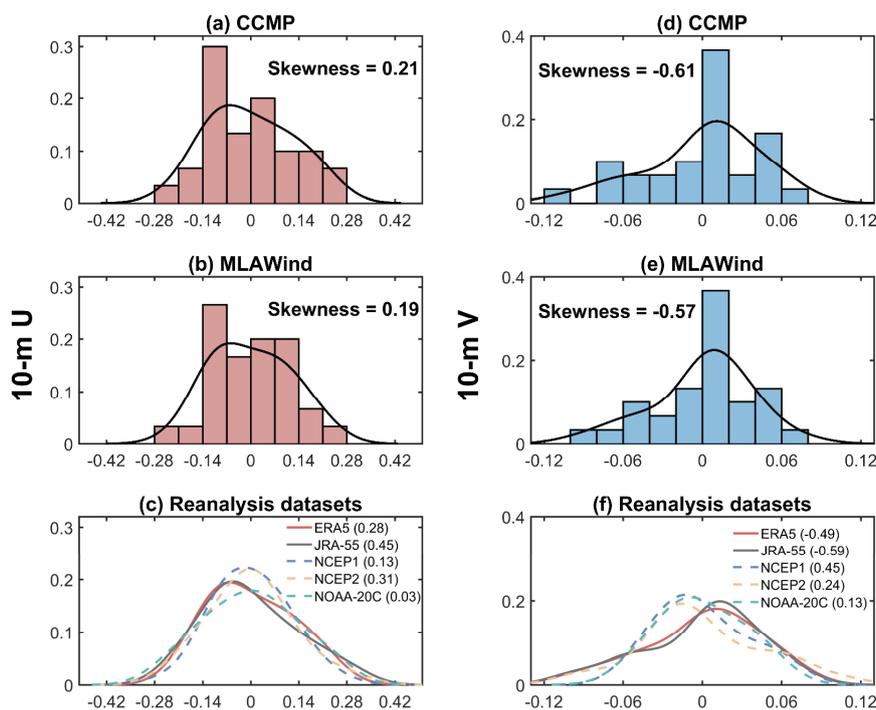


286 Atlantic Oceans. Their performance is comparable to those of near-global  
 287 coverage for both 10-m zonal (Fig. 4b-d) and meridional wind (Fig. 4f-h).



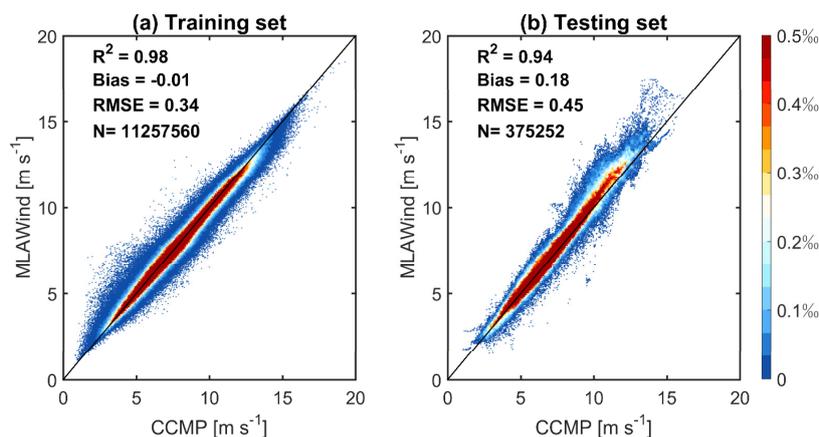
288  
 289 **Figure 4.** Evaluation based on satellite sea surface wind. (a-d) Bias and RMSE values ( $\text{m s}^{-1}$ ) of  
 290 multiple datasets evaluated with CCMP 10-m zonal wind over different oceans within  $60^{\circ}\text{S}$ – $60^{\circ}\text{N}$   
 291 from 1993–2022. (e-h) Same as (a-d) but for 10-m meridional wind.

292 The probability density function (PDF) characteristics of different sea  
 293 surface wind datasets are compared from 1993–2022. The 10-m zonal wind in  
 294 CCMP, MLAWind, and reanalysis datasets collectively exhibit positive  
 295 skewness distributions with longer tails toward larger values and left-shifted  
 296 peaks (Fig. 5a-c). MLAWind demonstrates the closest skewness value (0.19) to  
 297 CCMP satellite data (0.21) among all datasets. The 10-m meridional wind in  
 298 MLAWind and CCMP are consistent in their negative skewness distributions,  
 299 whereas substantial divergences emerge across reanalysis datasets (Fig. 5d-f).  
 300 In contrast to the CCMP dataset, the 10-m meridional wind in NCEP1, NCEP2,  
 301 and NOAA-20C exhibit positive skewness in their distributions (Fig. 5f), which  
 302 may be attributed to their large RMSE with satellite data (Fig. 4).



303  
304 **Figure 5.** Comparison of probability density functions (PDF) among different datasets. (a-c) PDF  
305 and skewness coefficients of 10-m zonal wind in different datasets from 1993–2022. The horizontal  
306 axis denotes 10-m zonal wind anomalies. The vertical axis denotes probability densities. The values  
307 in parentheses in (c) denote skewness coefficients of reanalysis datasets. (d-f) Same as (a-c) but for  
308 10-m meridional wind.

309 During the testing phase, the MLAWind dataset maintains strong  
310 agreement with satellite observations in 2023 (Fig. 6b). Most values fall near  
311 the 1:1 reference line, with relatively larger errors under high wind conditions  
312 ( $\geq 15$  m/s). The mean bias and RMSE between the MLAWind and CCMP  
313 datasets are  $0.18 \text{ m s}^{-1}$  and  $0.45 \text{ m s}^{-1}$ , respectively. Notably, the MLAWind  
314 dataset achieves R-squared values exceeding 0.90 for both the training and  
315 testing sets (Fig. 6a-b). The high consistency exhibits its robust generalization  
316 capabilities without signs of overfitting.



317  
318 **Figure 6.** Validation of MLAWind's generalization capability. Density scatter plots compare  
319 MLAWind and CCMP data during (a) training phase (1993–2022) and (b) testing phase (2023),  
320 respectively. The shading represents data density within each bin.  $R^2$ , mean bias, RMSE and sample  
321 size (N) are used as quantitative metrics for evaluating MLAWind's performance.

## 322 **4. Performance of MLAWind from climatology to long-term** 323 **trend**

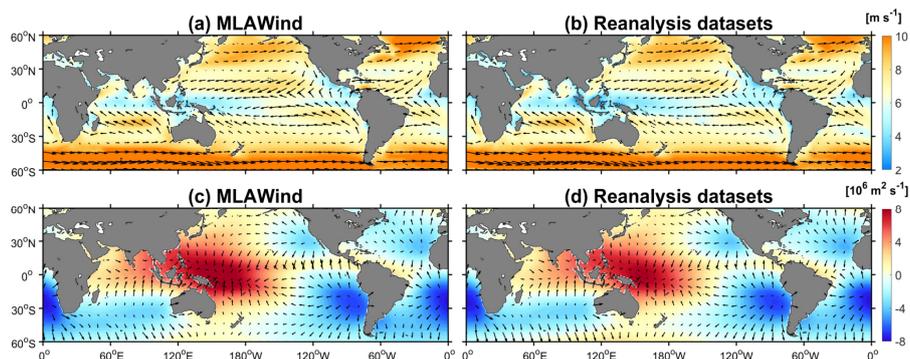
324 The reliability of the MLAWind dataset is validated through satellite  
325 observations in both training and non-training periods (Figs. 4-6). In this  
326 section, the performance of the MLAWind dataset since 1950 are assessed  
327 across different time scales, including climatology, seasonality, inter-annual  
328 variations, and long-term trends. Its similarities and differences with reanalysis  
329 datasets are evaluated, particularly during the period without satellite data.

### 330 **4.1 Climatology and annual cycle**

331 The climatological spatial pattern of the MLAWind dataset aligns closely  
332 with those of multiple reanalysis products, demonstrating robust consistency in  
333 capturing the tropical trade wind belt and the mid-latitude westerly belt (Fig.  
334 7a, b). Moreover, it can reproduce the key dynamical features of atmospheric  
335 circulation. The divergent wind and velocity potential of the MLAWind dataset  
336 exhibit a well-defined convergence–divergence dipole structure of the Walker  
337 circulation: intense convergence is concentrated over the Indo-Pacific Warm



338 Pool, while robust divergence centers are distinctly localized in the  
339 southeastern Pacific Ocean (Fig. 7c).

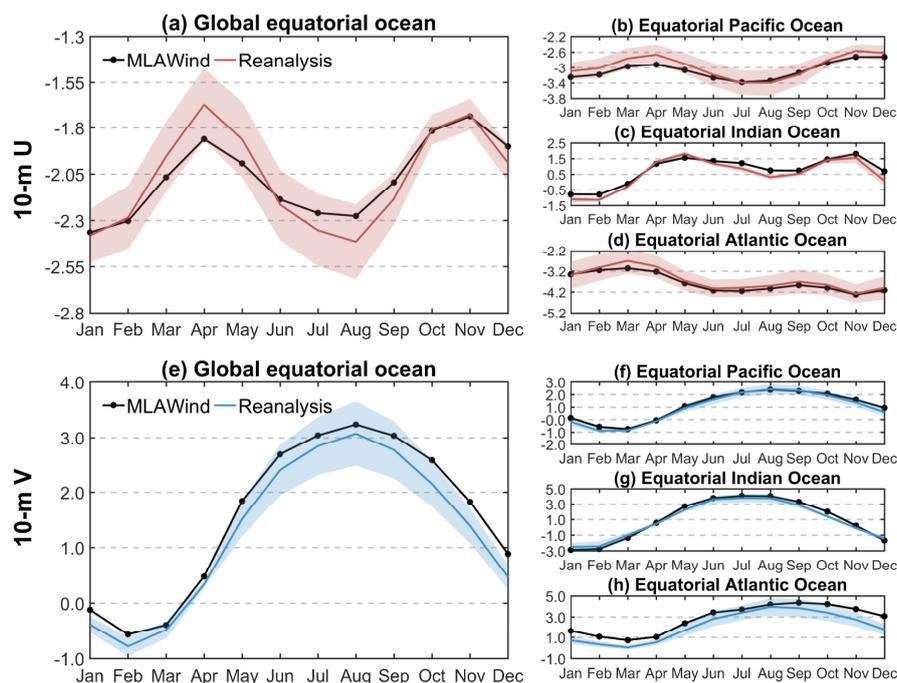


340  
341 **Figure 7.** (a-b) Climatology of 10-m wind velocity (vector,  $\text{m s}^{-1}$ ) and sea surface wind speed  
342 (shading,  $\text{m s}^{-1}$ ) in different datasets from 1981 to 2010. (c-d) Same as (a-b) but for divergent wind  
343 (vector,  $\text{m s}^{-1}$ ) and velocity potential (shading,  $\text{m}^2 \text{s}^{-1}$ ). The 30-year period is selected because all  
344 datasets are available during this period (Table 1). The climatological distribution of the MLAWind  
345 dataset over the 73-year period from 1950 to 2022 are examined and show no significant differences.

346 The annual cycle of mean sea surface winds is evaluated in equatorial  
347 regions within  $5^{\circ}\text{S}$ – $5^{\circ}\text{N}$  (Fig. 8). The 10-m zonal winds in the MLAWind  
348 dataset and reanalysis datasets show similar seasonal variations over the global  
349 equatorial ocean, with their easterlies persisted within a range of  $1.55$ – $2.55 \text{ m}$   
350  $\text{s}^{-1}$  throughout the year. A discrepancy in intensity exists between the two  
351 datasets. The 10-m zonal winds in the MLAWind dataset exhibit stronger  
352 easterlies during March–May and weaker easterlies during July–September  
353 compared to reanalysis datasets (Fig. 8a). The annual cycle characteristics of  
354 the MLAWind dataset show a general agreement with those of reanalysis  
355 datasets over the three equatorial oceans, albeit with some differences in  
356 magnitude (Fig. 8b-d). The 10-m meridional wind in both the MLAWind  
357 dataset and reanalysis datasets display a distinct unimodal structure, featuring  
358 a southerly wind peak in February and a northerly wind peak in August (Fig.  
359 8e). The MLAWind dataset accurately captures the annual cycle of 10-m  
360 meridional winds across different equatorial oceans, consistent with reanalysis



361 datasets (Fig. 8f-h).



362

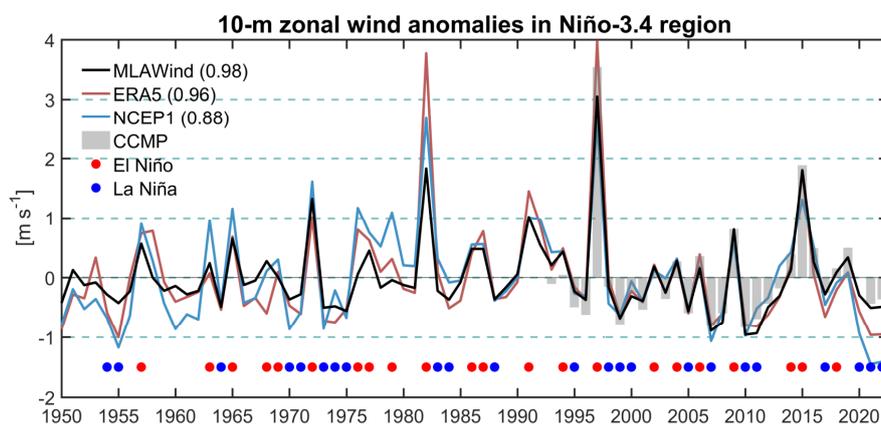
363 **Figure 8.** Annual cycle of sea surface wind ( $\text{m s}^{-1}$ ) over equatorial oceans ( $5^{\circ}\text{S}$ – $5^{\circ}\text{N}$ ) from 1981 to  
364 2010: (a-d) 10-m zonal wind and (e-h) 10-m meridional wind. The red shadings in (a-d) and the  
365 blue shading in (e-h) denote one standard deviations of 10-m zonal and meridional winds from  
366 multiple reanalysis datasets, respectively. The 30-year period (1981–2010) is selected because all  
367 datasets are available during this period (Table 1). The annual cycle of the MLAWind dataset over  
368 the 73-year period from 1950 to 2022 are examined and show no significant differences.

#### 369 4.2 Inter-annual variabilities

370 ENSO is the dominant inter-annual climate variability in the tropics,  
371 exerting significant influences on both regional and global climate systems (e.g.  
372 Bjerknes, 1969; Trenberth et al., 1998; McPhaden et al., 2006). Sea surface  
373 wind anomalies are greatly associated with the ENSO cycle (Bjerknes, 1969;  
374 McPhaden et al., 2006; Kug et al., 2009). Time series of 10-m zonal wind  
375 anomalies from multiple datasets are compared over the Niño-3.4 region (Fig.  
376 9). During the satellite observation period (1993–2022), MLAWind, ERA5, and



377 NCEP1 show strong agreement with CCMP, among which MLAWind exhibits  
378 the highest correlation coefficient (CORR = 0.98). During the historical period  
379 lacking satellite data (1950–1992), MLAWind maintains comparable inter-  
380 annual variations to those of both ERA5 and NCEP1, while demonstrating a  
381 higher correlation coefficient with the Niño-3.4 index (CORR = -0.86). The 10-  
382 m zonal wind anomalies derived from MLAWind show a robust association  
383 with ENSO dynamics, exhibiting strong westerly wind anomalies during El  
384 Niño events and persistent easterly anomalies during La Niña events, which  
385 aligns closely with the canonical Bjerknes feedback mechanism (Bjerknes,  
386 1969; McPhaden et al., 2006; Kug et al., 2009). Moreover, the inter-annual  
387 variabilities of MLAWind across other tropical ocean basins have been  
388 validated, equally exhibiting high consistency with reanalysis datasets and SST  
389 variabilities (Guo et al., 2025a).

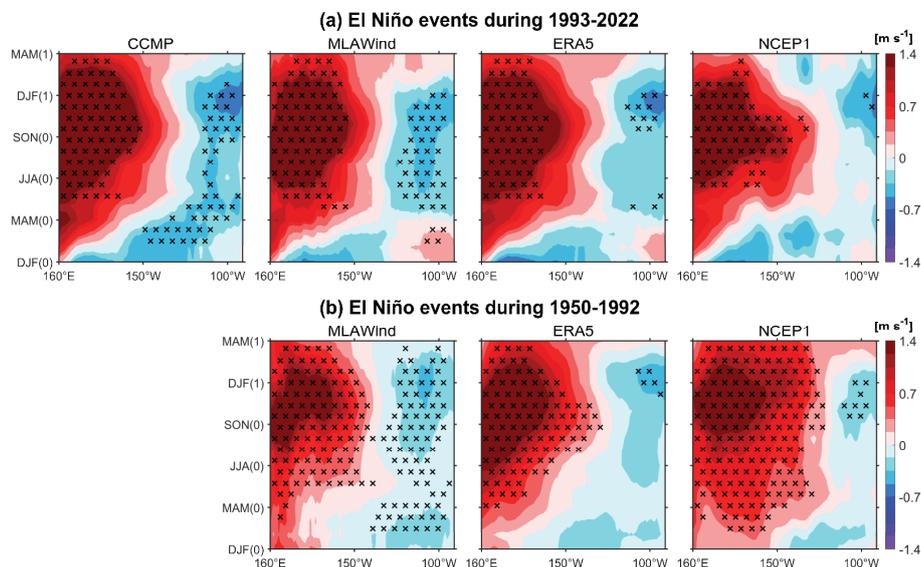


390  
391 **Figure 9.** Inter-annual variations of the November–December mean 10-m zonal wind anomalies ( $\text{m s}^{-1}$ ) over the Niño-3.4 region ( $5^{\circ}\text{S}$ – $5^{\circ}\text{N}$ ,  $170^{\circ}$ – $120^{\circ}\text{W}$ ). Long-term linear trends from 1950–2022  
392  $\text{s}^{-1}$ ) over the Niño-3.4 region ( $5^{\circ}\text{S}$ – $5^{\circ}\text{N}$ ,  $170^{\circ}$ – $120^{\circ}\text{W}$ ). Long-term linear trends from 1950–2022  
393 are removed. The grey bars denote 10-m zonal wind anomalies of CCMP from 1993–2022. The  
394 values in parentheses denote correlation coefficients with the CCMP data. The red and blue circles  
395 denote El Niño and La Niña events, respectively. El Niño and La Niña events are defined when the  
396 November–December mean Niño-3.4 index  $\geq +0.5^{\circ}\text{C}$  and  $\leq -0.5^{\circ}\text{C}$ , respectively.

397 A composite analysis of El Niño events is conducted for the satellite (Fig.



398 10a) and non-satellite periods (Fig. 10b). The selected El Niño events are  
399 marked by red circles in Fig. 9, with a total of 13 events during 1950–1992 and  
400 8 events during 1993–2022. During the satellite period (1993–2022),  
401 MLAWind reveal the development of westerly wind anomalies over the  
402 equatorial central-western Pacific Ocean ( $5^{\circ}\text{S}$ – $5^{\circ}\text{N}$ ,  $160^{\circ}\text{E}$ – $150^{\circ}\text{W}$ ), peaking  
403 around DJF(1). Weak easterly wind anomalies are observed over the equatorial  
404 eastern Pacific Ocean ( $5^{\circ}\text{S}$ – $5^{\circ}\text{N}$ ,  $150^{\circ}$ – $90^{\circ}\text{W}$ ). The pattern is consistent with  
405 those of CCMP, with the spatial correlation coefficient exceeding 0.95. ERA5  
406 and NCEP1 exhibit pronounced westerly anomalies in the region west of  $150^{\circ}\text{W}$ ,  
407 but the easterly anomalies in the equatorial eastern Pacific Ocean are weaker  
408 than those of CCMP and MLAWind (Fig. 10a).



409  
410 **Figure 10.** Composites of meridional mean ( $5^{\circ}\text{S}$ – $5^{\circ}\text{N}$ ) of 10-m zonal wind anomalies. (a) El Niño  
411 events during 1993–2022. (b) Same as (a) but for the period lacking satellite data (1950–1992).  
412 DJF(0)–SON(0) denote the periods in the developing year while DJF(1)–MAM(1) denote the  
413 periods in the decaying year. The black crosses denote the composites exceeding the 99%  
414 significance level based on Student's *t*-test.

415 During the historical period of 1950–1992, both MLAWind and reanalysis  
416 datasets show comparable sea surface wind patterns to those observed during



417 the satellite period. They collectively demonstrate significant westerly wind  
418 anomalies that are strongly linked to the development of El Niño events (Fig.  
419 10b). It is noted that there are several discrepancies in the three datasets.  
420 MLAWind shows relatively stronger easterly anomalies in the eastern  
421 equatorial Pacific Ocean. NCEP1 exhibits more pronounced westerly anomalies  
422 during the DJF(0)-MAM(0) period, which extends further eastward compared  
423 to those of MLAWind and ERA5 (Fig. 10b). Although the absence of robust  
424 benchmarks hinders the assessment of their relative accuracy, the consistency  
425 with satellite data and the Niño-3.4 index (Fig. 9) suggests that MLAWind may  
426 provide more reliable sea surface wind variations.

#### 427 **4.3 Long-term trends**

428 The long-term trend can filter the interference from short-term  
429 fluctuations, and thus reveal the linear change over a long time span (usually  
430 decades or longer periods). The long-term trends of the MLAWind dataset are  
431 evaluated in both the satellite (after 1993) and non-satellite (before 1993)  
432 periods to analyze its reliability across the two periods. The analysis during  
433 1993–2022 reveals that the MLAWind dataset attains the most robust spatial  
434 concordance with observations among all evaluated datasets (CORR = 0.95).  
435 This strong agreement is due to its low bias and RMSE with satellite data (Fig.  
436 4).

437 The long-term trends among different products exhibit significant  
438 divergence during historical periods without satellite data (Zhang et al., 2023;  
439 Guo et al., 2025a). Coral records are used to assess the performance of sea  
440 surface wind during these periods. The skeletal Mn/Ca ratio in corals from  
441 atolls with west-facing lagoons serves as an effective proxy for identifying  
442 westerly wind anomalies (Shen et al., 1992), thereby reflecting the strength and  
443 variability of zonal winds at both inter-annual and multi-decadal time scales  
444 (Thompson et al., 2015). The employed coral Mn/Ca records are collected from



445 three distinct sites in the tropical Pacific Ocean, including Tarawa Atoll  
446 (1.33°N, 172.97°E), Butaritari Atoll (3.07°N, 172.75°E) and Kiritimati Island  
447 (1.93°N, 157.49°W).

448 The coral Mn/Ca records from Tarawa Atoll (4°S–6°N, 166°–176°E)  
449 provide a valuable proxy for the period of 1950–1982, when satellite  
450 observations are unavailable. The remarkable increasing trend of Mn/Ca  
451 indicates the strengthening of the westerlies from 1950–1982. However,  
452 significant uncertainties remain in the long-term trends across multiple datasets  
453 (Table 2). Consistent with Mn/Ca records, the MLAWind and ERA5 datasets  
454 exhibit significant intensification of westerly winds ( $p < 0.05$ ), with the trend  
455 in MLAWind being more pronounced. The NCEP1 and NOAA-20C datasets  
456 show weaker trends with no statistical significance, demonstrating their  
457 relatively lower reliability for the period lacking satellite data. The long-term  
458 trends of Mn/Ca records from Butaritari Atoll (1989–2010) and Kiritimati  
459 Island (1994–2011) are examined, and the MLAWind dataset exhibits strong  
460 agreement.

461 **Table 2.** Long-term trends and the corresponding  $p$  values of coral Mn/Ca and 10-m zonal wind  
462 surrounding Tarawa Atoll (4°S–6°N, 166°–176°E) from 1950 to 1982. Positive values of coral  
463 Mn/Ca indicate an intensification of westerly winds.

Variable	Long-term trend	$P$ value
Coral record (nmol mol <sup>-1</sup> /decade)	<b>Mn/Ca</b>	<b>1.89</b>
	<b>MLAWind</b>	<b>0.49</b>
10-m U (m s <sup>-1</sup> /decade)	<b>ERA5</b>	<b>0.28</b>
	NCEP1	0.14
	NOAA-20C	0.14

464 A comprehensive analysis of sea surface wind trends in the tropical Indian  
465 Ocean (TIO) reveals significant discrepancies among multiple datasets (Guo et  
466 al., 2025a). Different from reanalysis results, the MLAWind dataset exhibits a



467 significant weakening of sea surface wind speed over the TIO during 1950–  
468 2019. The trend is corroborated by the asymmetrical west-east variations in  
469 both thermocline depth and sea surface height since 1950 (Guo et al., 2025a),  
470 demonstrating the reliability of the MLAWind dataset in capturing sea surface  
471 wind variations over long time periods.

## 472 **5. Data availability**

473 The MLAWind dataset reconstructed in this study is available at  
474 <https://doi.org/10.5281/zenodo.17354864> (Guo et al., 2025b). Here, we provide  
475 a monthly gridded sea surface wind product at  $1^\circ \times 1^\circ$  horizontal resolution from  
476 1950 to 2023, covering the near-global ocean between  $60^\circ\text{S}$  and  $60^\circ\text{N}$ .

## 477 **6. Conclusions**

478 This study employs a random forest algorithm integrated with the SHAP  
479 interpretable module to generate a monthly sea surface wind dataset  
480 (MLAWind). The MLAWind dataset effectively combines the respective  
481 strengths of in-situ and remote sensing observations, covering the near-global  
482 ocean between  $60^\circ\text{S}$  and  $60^\circ\text{N}$  from 1950 to 2023, with a spatial resolution of  
483  $1^\circ \times 1^\circ$ . The evaluations based on satellite data reveal that MLAWind exhibits  
484 better agreement with observations during the training period (1993–2022)  
485 compared to existing reanalysis datasets. It shows strong generalization  
486 capability, maintaining robust performance during the testing period (2023).  
487 The performance of MLAWind since 1950 is evaluated across different time  
488 scales, including climatology, annual cycle, inter-annual variation, and long-  
489 term trends. The assessments based on reanalysis datasets and coral records  
490 indicate that MLAWind demonstrates better skill in capturing long-term  
491 variations of sea surface winds compared to other datasets. MLAWind is  
492 expected to become an important resource for global climate change research.

493 The MLAWind data will continue to be optimized and updated in future  
494 work. The current version of the MLAWind dataset is limited within  $60^\circ\text{S}$ –



495 60°N, owing to the insufficient data coverage of the ICOADS and CCMP  
496 datasets in high-latitude regions. More reliable high-latitude datasets will be  
497 integrated into the machine learning model in the future, enabling the  
498 development of a global sea surface wind dataset. In addition, the MLAWind  
499 dataset, currently available at a horizontal resolution of  $1^\circ \times 1^\circ$ , demonstrates  
500 robust capability in capturing large-scale circulation signals. Further  
501 enhancement of its spatial resolution to  $0.25^\circ \times 0.25^\circ$  will significantly advance  
502 mesoscale marine and climate research. Several potential approaches can be  
503 used to construct a higher-resolution version of the MLAWind dataset,  
504 including the integration of higher-resolution input variables (Tian et al., 2022),  
505 the application of machine learning-based interpolation (He et al., 2022) and  
506 super-resolution techniques (Wang et al., 2021; Doury et al., 2023). An  
507 assessment of diverse approaches will be conducted in future work to ensure  
508 the upcoming high-resolution MLAWind dataset retains the existing version's  
509 accuracy while capturing mesoscale atmospheric processes.

#### 510 **Author contributions**

511 WG developed the data reconstruction model, generated the MLAWind dataset  
512 and drafted the initial manuscript. RZ participated in data analysis and revised the  
513 manuscript. XW conceptualized this study and contributed to the interpretation of the  
514 results. DW provided guidance on improving the manuscript. All authors provided  
515 valuable feedback and helped shape this study.

#### 516 **Competing interests**

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

#### 518 **Acknowledgements**

519 The authors would like to express their gratitude for the publicly available  
520 data resources that are utilized in this study. ICOADS is obtained from [https://  
521 /icoads.noaa.gov/products.html](https://icoads.noaa.gov/products.html). CCMP is obtained from <https://data.remss.com/cc>



522 [mp/](#). ERA5 is obtained from [https://cds.climate.copernicus.eu/datasets/reanalysis-e](https://cds.climate.copernicus.eu/datasets/reanalysis-era5-pressure-levels-monthly-means?tab=overview)  
523 [ra5-pressure-levels-monthly-means?tab=overview](#). NCEP1 is obtained from [https://](https://psl.noaa.gov/data/gridded/data.ncep.reanalysis.html)  
524 [psl.noaa.gov/data/gridded/data.ncep.reanalysis.html](#). NCEP2 is obtained from [http](https://psl.noaa.gov/data/gridded/data.ncep.reanalysis2.html)  
525 [s://psl.noaa.gov/data/gridded/data.ncep.reanalysis2.html](#). JRA-55 is obtained from  
526 <https://climatedataguide.ucar.edu/climate-data/jra-55>. NOAA-20C is obtained from  
527 [https://psl.noaa.gov/data/gridded/data.20thC\\_ReanV2c.html](https://psl.noaa.gov/data/gridded/data.20thC_ReanV2c.html). Coral records is obtai  
528 ned from <https://www.ncei.noaa.gov/products/paleoclimatology/coral-sclerosponge>.

### 529 **Financial support**

530 This work is supported by the National Natural Science Foundation of China  
531 (Grant No. 42330404), the National Key R&D Program of China (2022YFF0801400),  
532 the National Natural Science Foundation of China (42376027, 42406197), and the  
533 Special Fund of South China Sea Institute of Oceanology of the Chinese Academy of  
534 Sciences (SCSIO2023QY01).

535



536 **References**

- 537 Bagnell, A. and DeVries, T.: 20(th) century cooling of the deep ocean  
538 contributed to delayed acceleration of Earth's energy imbalance, *Nat.*  
539 *Commun.*, 12, 4604, <https://doi.org/10.1038/s41467-021-24472-3>, 2021.
- 540 Berry, D. I. and Kent, E. C.: Air–Sea fluxes from ICOADS: the construction of  
541 a new gridded dataset with uncertainty estimates, *Int. J. Clim.*, 31, 987–  
542 1001, <https://doi.org/10.1002/joc.2059>, 2011.
- 543 Bjerknes, J.: Atmospheric teleconnections from the equatorial Pacific, *Mon.*  
544 *Weather. Rev.*, 97, 163–172, [https://doi.org/10.1175/1520-0493\(1969\)097<0163:ATFTEP>2.3.CO;2](https://doi.org/10.1175/1520-0493(1969)097<0163:ATFTEP>2.3.CO;2), 1969.
- 546 Breiman, L.: Random Forests, *Machine Learning*, 45, 5–32,  
547 <https://doi.org/10.1023/A:1010933404324>, 2001.
- 548 Clarke, A. J.: El Niño physics and El Niño predictability, *Ann. Rev. Mar. Sci.*,  
549 6, 79–99, <https://doi.org/10.1146/annurev-marine-010213-135026>, 2014.
- 550 Compo, G. P., Whitaker, J. S., Sardeshmukh, P. D., Matsui, N., Allan, R. J., Yin,  
551 X., Gleason, B. E., Vose, R. S., Rutledge, G., Bessemoulin, P., Brönnimann,  
552 S., Brunet, M., Crouthamel, R. I., Grant, A. N., Groisman, P. Y., Jones, P.  
553 D., Kruk, M. C., Kruger, A. C., Marshall, G. J., Mauerer, M., Mok, H. Y.,  
554 Nordli, Ø., Ross, T. F., Trigo, R. M., Wang, X. L., Woodruff, S. D., and  
555 Worley, S. J.: The Twentieth Century Reanalysis Project, *Q. J. R. Meteorol.*  
556 *Soc.*, 137, 1–28, <https://doi.org/10.1002/qj.776>, 2011.
- 557 Deser, C., Alexander, M. A., and Timlin, M. S.: Evidence for a wind-driven  
558 intensification of the Kuroshio Current extension from the 1970s to the  
559 1980s, *J. Clim.*, 12, 1697–1706, [https://doi.org/10.1175/1520-0442\(1999\)012<1697:EFAWDI>2.0.CO;2](https://doi.org/10.1175/1520-0442(1999)012<1697:EFAWDI>2.0.CO;2), 1999.
- 561 Doury, A., Somot, S., Gadat, S., Ribes, A., and Corre, L.: Regional climate  
562 model emulator based on deep learning: concept and first evaluation of a  
563 novel hybrid downscaling approach, *Clim. Dyn.*, 60, 1751–1779,  
564 <https://doi.org/10.1007/s00382-022-06343-9>, 2023.



- 565 Findell, K. L., Keys, P. W., Ent R. J. V. D., Lintner, B. R., Berg, A., and Krasting,  
566 J. P.: Rising temperatures increase importance of oceanic evaporation as a  
567 source for continental precipitation, *J. Clim.*, 32, 7713–7726,  
568 <https://doi.org/10.1175/JCLI-D-19-0145.1>, 2019.
- 569 Freeman, E., Woodruff, S. D., Worley, S. J., Lubker, S. J., Kent, E. C., Lubker,  
570 W. E., Berry, D. I., Brohan, P., Eastman, R., Gates, L., Gloeden, W., Ji, Z.,  
571 Lawrimore, J., Rayner, N. A., Rosenhagen, G., and Smith, S. R.: ICOADS  
572 Release 3.0: a major update to the historical marine climate record, *Int. J.*  
573 *Clim.*, 37, 2211–2232, <https://doi.org/10.1002/joc.4775>, 2017.
- 574 Garreaud, R., Lopez, P., Minvielle, M., and Rojas, M.: Large-scale cont  
575 rol on the Patagonian climate, *J. Clim.*, 26, 215–230, [https://doi.or](https://doi.org/10.1175/JCLI-D-12-00001.1)  
576 [g/10.1175/JCLI-D-12-00001.1](https://doi.org/10.1175/JCLI-D-12-00001.1), 2013.
- 577 Guo, W., Zhang, R., Wang, X., Wang, C., Li, X., Han, W., and Zhang, L.:  
578 Unveiling the drivers of tropical Indian Ocean warming through machine  
579 learning-assisted surface wind, *J. Clim.*, 38, 6763–6779,  
580 <https://doi.org/10.1175/JCLI-D-25-0003.1>, 2025a.
- 581 Guo, W., Zhang, R., and Wang, X.: Machine learning-assisted sea surfa  
582 ce wind dataset (MLAWind), Zenodo [data set], [https://doi.org/10.52](https://doi.org/10.5281/zenodo.17354864)  
583 [81/zenodo.17354864](https://doi.org/10.5281/zenodo.17354864), 2025b.
- 584 He, Q., Wang, M., Liu, K., Li, K., and Jiang, Z.: GPRChinaTemp1km: a high-  
585 resolution monthly air temperature dataset for China (1951–2020) based  
586 on machine learning, *Earth Syst. Sci. Data.*, 14, 3273–3292,  
587 <https://doi.org/10.5194/essd-2021-442>, 2022.
- 588 Held, I. M. and Soden, B. J.: Robust responses of the hydrological cyc  
589 le to global warming, *J. Clim.*, 19, 5686–5699, [https://doi.org/10.11](https://doi.org/10.1175/JCLI3990.1)  
590 [75/JCLI3990.1](https://doi.org/10.1175/JCLI3990.1), 2006.
- 591 Hersbach, H., Bell, B., Berrisford, P., Horányi, A., Muñoz Sabater, J., Nicolas,  
592 J., Radu, R., Schepers, D., Simmons, A., Soci, C., and Dee, D.: Global  
593 reanalysis: goodbye ERA-Interim, hello ERA5, *Meteorology section of*



- 594 ECMWF Newsletter, 159, 17–24, <https://doi.org/10.21957/vf291hehd7>,  
595 2019.
- 596 Huang, B., Yin, X., Boyer, T., Liu, C., Menne, M., Rao, Y. D., Smith, T., Vose,  
597 R., and Zhang, H.-M.: Extended Reconstructed Sea Surface Temperature  
598 Version 6 (ERSSTv6): Part I. An Artificial Neural Network Approach, *J.*  
599 *Clim.*, 38, 1105–1121, <https://doi.org/10.1175/JCLI-D-23-0707.1>, 2025.
- 600 IPCC: Climate Change 2021: The Physical Science Basis, in: Contribution of  
601 Working Group I to the Sixth Assessment Report of the Intergovernmental  
602 Panel on Climate Change, edited by: Masson-Delmotte, V., Zhai, P., Pirani,  
603 A., Connors, S. L., Péan, C., Berger, S., Caud, N., Chen, Y., Goldfarb, L.,  
604 Gomis, M. I., Huang, M., Leitzell, K., Lonnoy, E., Matthews, J. B. R.,  
605 Maycock, T. K., Waterfield, T., Yelekçi, O., Yu, R., and Zhou, B.,  
606 Cambridge University Press, 2021.
- 607 Jiang, S., Sweet, L.-b., Blougouras, G., Brenning, A., Li, W., Reichstein, M.,  
608 Denzler, J., Shangguan, W. Yu, G., Huang, F., and Zscheischler, J.: How  
609 interpretable machine learning can benefit process understanding in the  
610 geosciences, *Earth's Future*, 12, e2024EF004540,  
611 <https://doi.org/10.1029/2024EF004540>, 2024.
- 612 Kalnay, E., Kanamitsu, M., Kistler, R., Collins, W., Deaven, D., Gandin, L.,  
613 Iredell, M., Saha, S., White, G., Woollen, J., Zhu, Y., Chelliah, M.,  
614 Ebisuzaki, W., Higgins, W., Janowiak, J., Mo, K. C., Ropelewski, C., Wang,  
615 J., Leetmaa, A., Reynolds, R., Jenne, R., and Joseph, D.: The NCEP/NCAR  
616 40-Year Reanalysis Project, *Bull. Am. Meteorol. Soc.*, 77, 437–471,  
617 [https://doi.org/10.1175/1520-0477\(1996\)077<0437:TNYRP>2.0.CO;2](https://doi.org/10.1175/1520-0477(1996)077<0437:TNYRP>2.0.CO;2),  
618 1996.
- 619 Kanamitsu, M., Ebisuzaki, W., Woollen, J., Yang, S.-K., Hnilo, J. J., Fiorino,  
620 M., and Potter, G. L.: NCEP–DOE AMIP-II Reanalysis (R-2), *Bull. Amer.*  
621 *Meteor. Soc.*, 83, 1631–1643, <https://doi.org/10.1175/BAMS-83-11-1631>,  
622 2002.



- 623 Kobayashi, S., Ota, Y., Harada, Y., Ebita, A., Moriya, M., Onoda, H., Onogi,  
624 K., Kamahori, H., Kobayashi, C., Endo, H., Miyaoka, K., and Takahashi,  
625 K.: The JRA-55 Reanalysis: General Specifications and Basic  
626 Characteristics, *J. Meteorol. Soc. Jpn. Ser. II*, 93, 5–48.  
627 <https://doi.org/10.2151/jmsj.2015-001>, 2015.
- 628 Kug, J.-S., Jin, F.-F., and An, S.-I.: Two types of El Niño events: Cold tongue  
629 El Niño and warm pool El Niño, *J. Clim.*, 22, 1499–1515,  
630 <https://doi.org/10.1175/2008JCLI2624.1>, 2009.
- 631 Leng, J., Chen, J. M., Li, W., Luo, X., Xu, M., Liu, J., Wang, R. Rogers, C. Li,  
632 B., and Yan, Y.: Global datasets of hourly carbon and water fluxes  
633 simulated using a satellite-based process model with dynamic  
634 parameterizations, *Earth Syst. Sci. Data*, 16, 1283–1300,  
635 <https://doi.org/10.5194/essd-16-1283-2024>, 2024.
- 636 Li, R., Li, Y., Lyu, Y., Sprintall, J., and Wang, F.: Role of the Indian Ocean  
637 Wind-Driven Dynamics in the Indonesian Throughflow Variability, *J.*  
638 *Geophys. Res.: Oceans*, 130, e2025JC022503,  
639 <https://doi.org/10.1029/2025JC022503>, 2025.
- 640 Lundberg, S. and Lee, S.-I.: A unified approach to interpreting model  
641 predictions, *Adv. Neural Inf. Process. Syst.*, 30, 1–10,  
642 <https://doi.org/10.48550/arXiv.1705.07874>, 2017.
- 643 Mao, W., Mu, X., Zheng, Y., and Yan, G.: Leave-one-out cross-validation-based  
644 model selection for multi-input multi-output support vector machine,  
645 *Neural Comput. App.*, 24, 441–451, [https://doi.org/10.1007/s00521-012-](https://doi.org/10.1007/s00521-012-1234-5)  
646 1234-5, 2014.
- 647 McPhaden, M. J., Zebiak, S. E., and Glantz, M. H.: ENSO as an integrating  
648 concept in Earth science, *Science*, 314, 1740–1745,  
649 <https://doi.org/10.1126/science.1132588>, 2006.
- 650 Mears, C. A., Scott, J., Wentz, F. J., Ricciardulli, L., Leidner, S. M., Hoffman,  
651 R., and Atlas, R.: A near-real-time version of the cross-calibrated



- 652 multiplatform (CCMP) ocean surface wind velocity data set, *J. Geophys.*  
653 *Res.: Oceans*, 124, 6997–7010, <https://doi.org/10.1029/2019JC015367>,  
654 2019.
- 655 Qin, L., Zhu, L., Liu, B., Li, Z., Tian, Y., Mitchell, G., Shen, S., Xu, W., and  
656 Chen, J.: Global expansion of tropical cyclone precipitation footprint, *Nat.*  
657 *Commun.*, 15, <https://doi.org/10.1038/s41467-024-49115-1>, 2024.
- 658 Reichstein, M., Camps-Valls, G., Stevens, B., Jung, M., Denzler, J., Carvalhais,  
659 N., and Prabhat: Deep learning and process understanding for data-driven  
660 Earth system science, *Nature*, 566, 195–204,  
661 <https://doi.org/10.1038/s41586-019-0912-1>, 2019.
- 662 Ren, H. L., Zuo, J., and Deng, Y.: Statistical predictability of Niño indices for  
663 two types of ENSO, *Clim. Dyn.*, 52, 5361–5382.  
664 <https://doi.org/10.1007/s00382-018-4453-3>, 2019.
- 665 Sayani, H. R., Thompson, D. M., Carilli, J. E., Marchitto, T. M., Chapman, A.  
666 U., and Cobb, K. M.: Reproducibility of coral Mn/Ca-based wind  
667 reconstructions at Kiritimati Island and Butaritari Atoll, *Geochem. Geophys.*  
668 *Geosy.*, 22, e2020GC009398, <https://doi.org/10.1029/2020GC009398>,  
669 2021.
- 670 Shen, G. T., Linn, L. J., Campbell, T. M., Cole, J. E., and Fairbanks, R. G.: A  
671 chemical indicator of trade wind reversal in corals from the western  
672 tropical Pacific, *J. Geophys. Res.: Oceans.*, 97, 12698–12697,  
673 <https://doi.org/10.1029/92JC00951>, 1992.
- 674 Su, H., Li, W. and Yan, X.-H.: Retrieving temperature anomaly in the global  
675 subsurface and deeper ocean from satellite observations, *J. Geophys. Res.:  
676 Oceans.*, 123, 399–410, <https://doi.org/10.1002/2017JC013631>, 2018.
- 677 Su, H., Zhang, H., Geng, X., Qin, T., Lu, W., and Yan, X. H.: OPEN: A new  
678 estimation of global ocean heat content for upper 2000 meters from remote  
679 sensing data, *Remote Sens.*, 12, 2294, <https://doi.org/10.3390/rs12142294>,  
680 2020.



- 681 Thompson, D. M., Cole, J. E., Shen, G. T., Tudhope, A. W., and Meehl, G. A.:  
682 Early twentieth-century warming linked to tropical Pacific wind strength,  
683 Nat. Geosci., 8, 117–121, <https://doi.org/10.1038/ngeo2321>, 2015.
- 684 Tian, T. Cheng, L., Wang, G., Abraham, J., Ren, S., Zhu, J., Song, J., and Leng,  
685 H.: Reconstructing ocean subsurface salinity at high resolution using a  
686 machine learning approach, Earth Syst. Sci. Data, 14, 5037–5060,  
687 <https://doi.org/10.5194/essd-2022-236>, 2022.
- 688 Tokinaga, H. and Xie, S.-P.: Wave- and Anemometer-Based Sea Surface Wind  
689 (WASWind) for Climate Change Analysis\*, J. Clim., 24, 267–285,  
690 <https://doi.org/10.1175/2010JCLI3789.1>, 2011.
- 691 Tokinaga, H., Xie, S.-P., Timmermann, A., McGregor, S., Ogata, T., Kubota, H.,  
692 and Okumura, Y. M.: Regional Patterns of Tropical Indo-Pacific Climate  
693 Change: Evidence of the Walker Circulation Weakening\*, J. Clim., 25,  
694 1689–1710, <https://doi.org/10.1175/JCLI-D-11-00263.1>, 2012.
- 695 Trenberth, K. E., Branstator, G. W., Karoly, D., Kumar, A., Lau, N. C.,  
696 and Ropelewski, C.: Progress during TOGA in understanding and  
697 modeling global teleconnections associated with tropical sea surface  
698 temperatures, J. Geophys. Res.: Oceans., 103, 14291–14324, <https://doi.org/10.1029/97jc01444>, 1998.
- 700 Tseng, Y. H., Huang, J. H., and Chen, H. C.: Improving the Predictability of  
701 Two Types of ENSO by the Characteristics of Extratropical Precursors,  
702 Geophys. Res. Lett., 49, e2021GL097190,  
703 <https://doi.org/10.1029/2021GL097190>, 2022.
- 704 Vinoth, J. and Young, I. R.: Global Estimates of Extreme Wind Speed and Wave  
705 Height, J. Clim., 24, 1647–1665, <https://doi.org/10.1175/2010JCLI3680.1>,  
706 2011.
- 707 Wang, C.: Three-ocean interactions and climate variability: A review and  
708 perspective, Clim. Dyn., 53, 18–5136, <https://doi.org/10.1007/s00382-019-04930-x>, 2019.



- 710 Wang, D., Liu, Q., Huang, R. X., Du, Y., and Qu T.: Interannual variability of  
711 the South China Sea throughflow inferred from wind data and an ocean  
712 data assimilation product, *Geophys. Res. Lett.*, 33, L14605,  
713 <https://doi.org/10.1029/2006GL026316>, 2006.
- 714 Wang, H. and Li, X.: DeepBlue: Advanced convolutional neural network  
715 k applications for ocean remote sensing. *IEEE Geoscience and Re  
716 mote Sensing Magazine*, 12, 138–161, [https://doi.org/10.1109/MGRS.  
717 2023.3343623](https://doi.org/10.1109/MGRS.2023.3343623), 2024.
- 718 Wang, J., Liu, Z., Foster, I., Chang, W., Kettimuthu, R., and Kotamarthi, V. R.:  
719 Fast and accurate learned multiresolution dynamical downscaling for  
720 precipitation, *Geosci. Model. Dev.*, 14, 6355–6372,  
721 <https://doi.org/10.5194/gmd-14-6355-2021>, 2021.
- 722 Wang, Q., Zeng, L., Chen, J., He, Y., Liu, Q., Sui, D., and Wang, D.: Phase shift  
723 of the winter South China Sea western boundary current over the past two  
724 decades and its drivers, *Geophys. Res. Lett.*, 50, e2023GL103145,  
725 <https://doi.org/10.1029/2023GL103145>, 2023.
- 726 Wang, X., Dickinson, R. E., Su, L., Zhou, C., and Wang, K.: PM2.5 pollution  
727 in China and how it has been exacerbated by terrain and meteorological  
728 conditions, *Bull. Am. Meteorol. Soc.*, 99, 105–119,  
729 <https://doi.org/10.1175/bams-d-16-0301.1>, 2018.
- 730 Watt-Meyer, O., Brenowitz, N. D., Clark, S. K., Henn, B., and Bretherton, C.  
731 S.: Correcting weather and climate models by machine learning nudged  
732 historical simulations, *Geophys. Res. Lett.*, 48, e2021GL092555,  
733 <https://doi.org/10.1029/2021GL092555>, 2021.
- 734 Yang, L., Wang, D., Huang, J., Wang, X., Zeng, L., Shi, R., He, Y., Xie, Q.,  
735 Wang, S., Chen, R., Yuan, J., Wang, Q., Chen, J., Zu, T., Li, J., Sui, D.,  
736 and Peng, S.: Toward a Mesoscale Hydrological and Marine  
737 Meteorological Observation Network in the South China Sea, *Bull. Amer.*



- 738 Meteor. Soc., 96, 1117–1135, <https://doi.org/10.1175/BAMS-D-14->  
739 00159.1, 2015.
- 740 Yin, J., Slater, L. J., Khouakhi, A., Yu, L., Liu, P., Li, F., Pokhrel, Y., and  
741 Gentine, P.: GTWS-MLrec: global terrestrial water storage reconstruction  
742 by machine learning from 1940 to present, *Earth Syst. Sci. Data*, 15, 5597–  
743 5615, <https://doi.org/10.5194/essd-15-5597-2023>, 2023.
- 744 Yoo, C., Im, J., Park, S., and Quackenbush, L. J.: Estimation of daily maximum  
745 and minimum air temperatures in urban landscapes using MODIS time  
746 series satellite data, *ISPRS J. Photogramm. Remote Sens.*, 137, 149–162,  
747 <https://doi.org/10.1016/j.isprsjprs.2018.01.018>, 2018.
- 748 Young, I. and Ribal, A.: Multiplatform evaluation of global trends in wind  
749 speed and wave height, *Science*, 364, 548,  
750 <https://doi.org/10.1126/science.aav9527>, 2019.
- 751 Zhang, H., Chen, D., Liu, T., Tian, D., He, M., Li, Q., Wei, G., and Liu, J.:  
752 MASCS 1.0: synchronous atmospheric and oceanic data from a cross-  
753 shaped moored array in the northern South China Sea during 2014–2015,  
754 *Earth Syst. Sci. Data*, 16, 5665–5679, <https://doi.org/10.5194/essd-16->  
755 5665-2024, 2024.
- 756 Zhang, R., Guo, W., Wang, X., and Wang, C.: Ambiguous Variations in Tropical  
757 Latent Heat Flux since the Years around 1998, *J. Clim.*, 36, 3403–3415,  
758 <https://doi.org/10.1175/JCLI-D-22-0381.1>, 2023.
- 759 Zhang, X., Nikurashin, M., Pena-Molino, B., Rintoul, S. R., and Doddridge, E.:  
760 A Theory of Standing Meanders of the Antarctic Circumpolar Current and  
761 Their Response to Wind, *J Phys. Oceanogr.*, 53, 235–251,  
762 <https://doi.org/10.1175/JPO-D-22-0086.1>, 2023.
- 763 Zhou, C., Azorin-Molina, C., Engstrm, E., Minola, L., Wern, L., Hellstrm, S.,  
764 Lnn, J., and Chen, D.: HomogWS-se: a century-long homogenized dataset  
765 of near-surface wind speed observations since 1925 rescued in Sweden,



766 Earth Syst. Sci. Data, 14, 2167–2177, <https://doi.org/10.5194/essd-14->  
767 2167-2022, 2022.