

1 **Assessing improvements in global ocean pCO<sub>2</sub> machine learning reconstructions with**  
2 **Southern Ocean autonomous sampling**

3 Thea H. Heimdal<sup>1</sup>, Galen A. McKinley<sup>1</sup>, Adrienne J. Sutton<sup>2</sup>, Amanda R. Fay<sup>1</sup>, Lucas Gloege<sup>3</sup>

4 <sup>1</sup>Columbia University and Lamont-Doherty Earth Observatory, Palisades, NY, USA

5 <sup>2</sup>Pacific Marine Environmental Laboratory, National Oceanic and Atmospheric Administration,  
6 Seattle, WA, USA

7 <sup>3</sup>Open Earth Foundation, Marina del Rey, CA, USA

8 *Correspondence to:* Thea H. Heimdal (theimdal@ldeo.columbia.edu)

9

10 **Abstract**

11 The Southern Ocean plays an important role in the exchange of carbon between the atmosphere  
12 and oceans, and is a critical region for the ocean uptake of anthropogenic CO<sub>2</sub>. However, estimates  
13 of the Southern Ocean air-sea CO<sub>2</sub> flux are highly uncertain due to limited data coverage. Increased  
14 sampling in winter and across meridional gradients in the Southern Ocean may improve machine  
15 learning (ML) reconstructions of global surface ocean pCO<sub>2</sub>. Here, we use a Large Ensemble  
16 Testbed (LET) of Earth System Models and the pCO<sub>2</sub>-Residual reconstruction method to assess  
17 improvements in pCO<sub>2</sub> reconstruction fidelity that could be achieved with additional autonomous  
18 sampling in the Southern Ocean added to existing Surface Ocean CO<sub>2</sub> Atlas (SOCAT)  
19 observations. The LET allows for a robust evaluation of the skill of pCO<sub>2</sub> reconstructions in space  
20 and time through comparison to ‘model truth’. With only SOCAT sampling, Southern Ocean and  
21 global pCO<sub>2</sub> are overestimated, and thus the ocean carbon sink is underestimated. Incorporating  
22 Uncrewed Surface Vehicle (USV) sampling increases the spatial and seasonal coverage of  
23 observations within the Southern Ocean, leading to a decrease in the overestimation of pCO<sub>2</sub>. A  
24 modest number of additional observations in southern hemisphere winter and across meridional  
25 gradients in the Southern Ocean leads to improvement in reconstruction bias and root-mean  
26 squared error (RMSE) by as much as 95 % and 16 %, respectively, as compared to SOCAT  
27 sampling alone. Lastly, the large decadal variability of air-sea CO<sub>2</sub> fluxes shown by SOCAT-only  
28 sampling may be partially attributable to undersampling of the Southern Ocean.

29

## 30 1. Introduction

31 The ocean plays an important role in mitigating climate change by sequestering anthropogenic  
32 carbon emissions. From 1850 to 2023, the oceans have removed a total of  $180 \pm 35$  Gt of carbon  
33 (Friedlingstein et al., 2023). In order to fully understand the climate impacts from rising emissions,  
34 it is essential to accurately quantify the air-sea  $\text{CO}_2$  flux and the global ocean carbon sink in space  
35 and time. The Surface Ocean  $\text{CO}_2$  Atlas (SOCAT; Bakker et al., 2016) is the largest global  
36 database of surface ocean  $\text{CO}_2$  observations, with data starting in 1957. The main synthesis and  
37 gridded products contain over 33 million high-quality direct shipboard measurements of  $f\text{CO}_2$   
38 (fugacity of  $\text{CO}_2$ ) with an uncertainty of  $< 5 \mu\text{atm}$  (Bakker et al., 2022). However, due to limited  
39 resources for ocean observing, limited number of ships/routes, inaccessible regions and unsafe  
40 waters, the database covers only about 1% of the global ocean at monthly  $1^\circ \times 1^\circ$  spatial resolution  
41 over the period of 1982-2023, and is highly biased towards the northern hemisphere.

42 Mapping methods have been developed to estimate full-coverage surface ocean  $p\text{CO}_2$   
43 across space and time by extrapolating to global coverage from these sparse SOCAT observations  
44 (e.g., Landschützer et al., 2014; Rödenbeck et al., 2015; Gloege et al., 2022; Bennington et al.,  
45 2022a,b). Most of these data products utilize machine learning (ML) algorithms to estimate a non-  
46 linear function between a suite of driver variables (i.e., sea surface temperature - SST, sea surface  
47 salinity - SSS, mixed layer depth - MLD, Chlorophyll - Chl-a,  $x\text{CO}_2$  - atmospheric  $\text{CO}_2$ ) and  
48 surface ocean  $p\text{CO}_2$  (the target variable) where these are co-located. The driver variables are  
49 proxies for processes influencing ocean  $p\text{CO}_2$ . Full-coverage driver variable datasets are then  
50 processed through these ML algorithms to produce estimated global full-coverage surface ocean  
51  $p\text{CO}_2$ . Since the data products rely on  $p\text{CO}_2$  observations to estimate functions between the target  
52 and driver variables, data sparsity remains a fundamental limitation to this technique.

53 It has been suggested that targeted sampling from autonomous platforms combined with  
54 ships, filling in the state space of  $p\text{CO}_2$ , represents a path forward to improve surface ocean  $p\text{CO}_2$   
55 reconstructions (Bushinsky et al., 2019; Gregor et al., 2019; Gloege et al., 2021; Djeutchouang et  
56 al., 2022; Landschützer et al., 2023; Hauck et al., 2023). One major obstacle, however, is that the  
57 indirect  $p\text{CO}_2$  estimates from floats have high uncertainties ( $\pm 11.4 \mu\text{atm}$ ) and may be biased by  
58 as much as  $\sim 4 \mu\text{atm}$  (Bakker et al., 2016; Williams et al., 2017; Fay et al., 2018; Gray et al., 2018;  
59 Sutton et al., 2021; Mackay and Watson 2021; Wu et al 2022). These large uncertainties and biases

60 arise when  $p\text{CO}_2$  is not measured directly as in the observations included in SOCAT, but is rather  
61 estimated using measurements of pH combined with a regression-derived alkalinity estimate  
62 (Williams et al., 2017; Gray et al., 2018). SOCAT includes only direct  $p\text{CO}_2$  observations. Biases  
63 and uncertainties may have large impacts on global air-sea  $\text{CO}_2$  flux estimates, given that the global  
64 mean air-sea disequilibrium is only 5-8  $\mu\text{atm}$  (McKinley et al., 2020). It is therefore critical that  
65 bias and uncertainty corrections are well-constrained over different oceanic conditions and over  
66 time.

67 Uncrewed Surface Vehicles (USVs), such as those manufactured and maintained by  
68 Saildrone Inc., represent a new type of autonomous platform that can obtain direct  $p\text{CO}_2$   
69 observations with significantly lower uncertainties compared to other autonomous methods, and  
70 equivalent to the highest-quality shipboard measurements contained in SOCAT ( $\pm 2 \mu\text{atm}$ ; Sabine  
71 et al., 2020; Sutton et al., 2021). Such improvements in sampling are critically important in the  
72 undersampled Southern Ocean. This region is fundamental in terms of the ocean's ability to  
73 remove carbon from the atmosphere, being responsible for  $\sim 40\%$  of the global ocean uptake of  
74 anthropogenic  $\text{CO}_2$  (Khatiwala et al., 2009). Improved data coverage in the Southern Ocean  
75 represents thus a major opportunity to advance our understanding of the global ocean carbon sink  
76 (Lenton et al., 2006, 2013; Takahashi et al., 2009; Monteiro et al., 2015; Gregor et al., 2019; Gray  
77 et al., 2018; Mongwe et al., 2018; Bushinsky et al., 2019; Sutton et al., 2021; Long et al., 2021;  
78 Mackay et al., 2022; Wu et al., 2022; Landschützer et al., 2023; Hauck et al., 2023). A combination  
79 of SOCAT and Saildrone USV observations would include high-accuracy data from both the long  
80 record and global coverage of ship tracks, and the expanded finer resolution of spatial and seasonal  
81 coverage of the poorly sampled Southern Ocean. Importantly, Saildrone USVs are also able to  
82 cover the spatial extent and seasonal cycle of the meridional gradients, which has been shown to  
83 be critical in order to reduce errors in reconstructing surface ocean  $p\text{CO}_2$  (Djeutchouang et al.,  
84 2022). A combined approach, with autonomous samples such as those obtained from Saildrone  
85 USVs, in addition to high-quality observations collected from ships, represents thus a promising  
86 solution to improve surface ocean  $p\text{CO}_2$  ML reconstructions.

87 Here, we assess to what extent surface ocean  $p\text{CO}_2$  reconstructions can improve by  
88 implementing the  $p\text{CO}_2$ -Residual machine learning (ML) reconstruction (Bennington et al., 2022a)  
89 with the combined inputs of SOCAT and Saildrone USV coverage. However, instead of using real-

90 world observations, we sample the target (i.e., surface ocean pCO<sub>2</sub>) and driver variables (i.e., SST,  
91 SSS, MLD, Chl-a and xCO<sub>2</sub>) from our Large Ensemble Testbed (LET) of Earth System Models  
92 (ESMs) (e.g., Stamell et al., 2020; Gloege et al., 2021; Bennington et al., 2022a). There are two  
93 major benefits of using a testbed compared to actual observations. First, in an ESM, the surface  
94 ocean pCO<sub>2</sub> field is provided precisely at all model times and 1°x1° points. Therefore, the pCO<sub>2</sub>  
95 reconstructed by the ML algorithm can be robustly evaluated in space and time against a known  
96 ‘truth’ (i.e., ‘model truth’). The reconstruction evaluation is thus not limited to the availability of  
97 sparse real-world ocean observations. Secondly, a testbed can be used to plan and evaluate the  
98 impact of different sampling strategies on the reconstructed pCO<sub>2</sub>. It is important to stress that, by  
99 using a model testbed, we do not predict real-world surface ocean pCO<sub>2</sub> and air-sea CO<sub>2</sub> fluxes.  
100 The goal here is to assess the accuracy with which an ML algorithm can reconstruct the ‘model  
101 truth’ given inputs of samples consistent with real-world data coverage from the SOCAT database  
102 and Saildrone USVs.

103 By utilizing the observational coverage of SOCAT and Saildrone USV transects, we assess  
104 to what extent the pCO<sub>2</sub>-Residual method accurately reconstructs model surface ocean pCO<sub>2</sub> in  
105 space and time. We test the impact of two different USV Southern Ocean sampling schemes, the  
106 first based on a sampling campaign completed in 2019 (Sutton et al., 2021), and the second on  
107 logistically feasible potential future meridional sampling. Additionally, we explore the timing,  
108 magnitude, duration and spatial extent of Southern Ocean USV sample additions that most  
109 significantly improve the pCO<sub>2</sub> predictions. Combined, the sampling patterns tested here  
110 complements previous studies exploring the impact of additional sampling in the Southern Ocean  
111 based on idealized full global coverage of floats, and float observations from recent deployments,  
112 including the Southern Ocean Carbon and Climate Observations and Modeling (SOCCOM)  
113 project, moorings and sailboats (Bushinsky et al., 2019; Denvil-Sommer et al., 2021;  
114 Djetchouang et al., 2022; Hauck et al., 2023; Behncke et al., 2024; Landschützer et al., 2023).

115

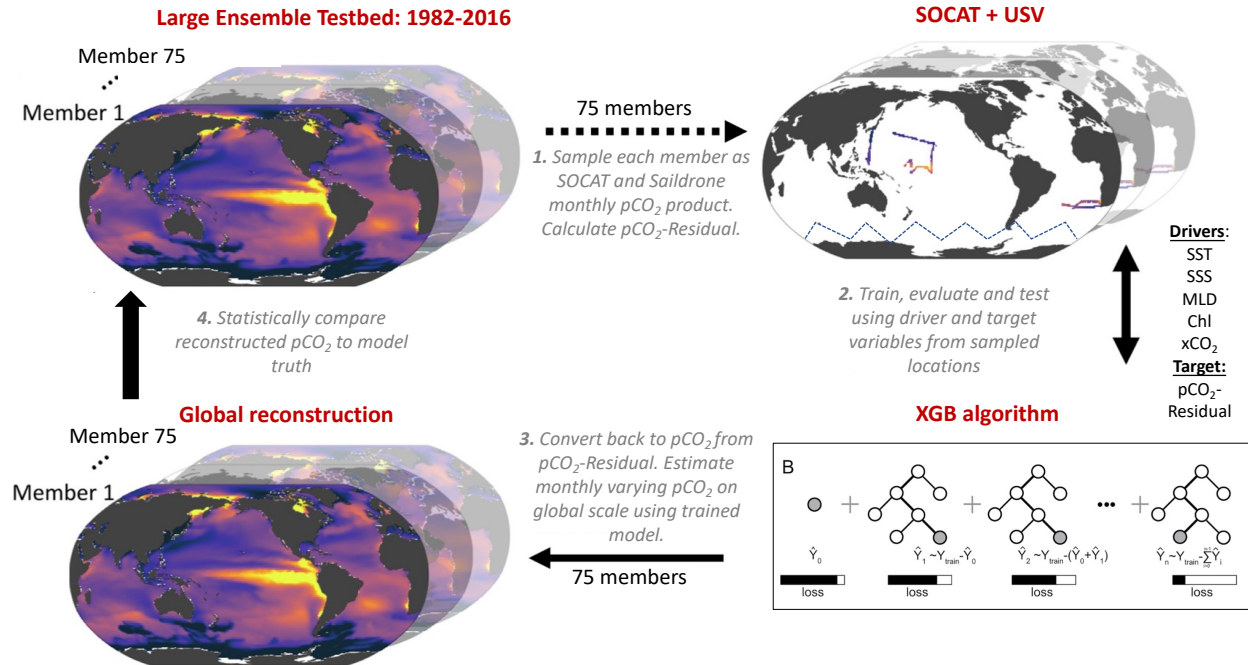
## 116 **2. Methods**

### 117 *2.1 The Large Ensemble Testbed (LET)*

118 In this study, the Large Ensemble Testbed (LET) includes 25 members from three independent  
119 initial-condition ensemble models (i.e., CanESM2, CESM-LENS and GFDL-ESM2M; Kay et al.,  
120 2015; Rodgers et al., 2015; Fyfe et al., 2017), giving a total of 75 members within the testbed. We  
121 do not use the MPI-GE model that was included in the past LET studies because its Southern  
122 Ocean pCO<sub>2</sub> seasonality and decadal variability appear to be anomalously large (Gloege et al.,  
123 2021; Fay and McKinley, 2021; Bennington et al., 2022a). Each individual Earth System Model  
124 (ESM) is an imperfect representation of the actual Earth system, so the multiple Large Ensembles  
125 are used to span different model structures and their representation of internal variability. Each  
126 ensemble member undergoes the same external forcing (i.e., historical atmospheric CO<sub>2</sub> before  
127 2005 and Representative Concentration Pathway 8.5 through 2016, plus solar and volcanic  
128 forcing), but the spread across the ensemble members gives a unique trajectory of the ocean-  
129 atmosphere state over time, i.e., a different state of internal variability as well as the difference  
130 across models.

131 The LET used in this study includes monthly 1°x1° model output from 1982-2016 (Gloege  
132 et al., 2021). For each individual ensemble member of the LET, surface ocean pCO<sub>2</sub> and co-located  
133 driver variables (i.e., SST, SSS, Chl-a, MLD, xCO<sub>2</sub>) were sampled monthly at a 1°x1° resolution,  
134 at times and locations equivalent to SOCAT and Saildrone USV observations (**Fig. 1**; Step 1).  
135 While the SOCAT observations were sampled from the testbed matching the actual years of  
136 sampling, the USV observations were sampled from the testbed starting in 2007 (for ten-year  
137 sampling) or 2012 (for five-year sampling) (see **Sect. 2.4**). As our focus is on reconstruction for  
138 the open ocean, testbed output for coastal areas, the Arctic Ocean (>79°N) and marginal seas  
139 (Hudson Bay, Caspian Sea, Black Sea, Mediterranean Sea, Baltic Sea, Java Sea, Red Sea and Sea  
140 of Okhotsk) were removed prior to algorithm processing.

141



142 **Figure 1:** Schematic of the Large Ensemble Testbed (LET; modified from Gloege et al., 2021). **1:** Surface ocean  
 143 pCO<sub>2</sub> from each of the 75 model members is sampled in space and time mimicking real-world SOCAT and Saildrone  
 144 USV observations (see **Fig. 2; Table 1; Section 2.5**). Prior to algorithm processing, pCO<sub>2</sub>-Residual is calculated  
 145 (**Section 2.2**). **2:** The pCO<sub>2</sub>-Residual (target variable) and co-located driver variables (i.e., SST, SSS, MLD, Chl,  
 146 xCO<sub>2</sub>) sampled from the testbed are processed by the XGBoost (XGB) algorithm (**Section 2.3**). **3:** Based on the full-  
 147 coverage of driver variables, pCO<sub>2</sub>-Residual is reconstructed globally. This process is repeated 75 times, individually  
 148 for every single testbed model member. The temperature component (pCO<sub>2</sub>-T) is then added back to the pCO<sub>2</sub>-  
 149 Residual for each value. **4:** The globally reconstructed pCO<sub>2</sub> is evaluated against the ‘model truth’ at all 1°x1° grid  
 150 cells. SST = sea surface temperature. SSS = sea surface salinity. MLD = mixed layer depth. Chl = chlorophyll. xCO<sub>2</sub>  
 151 = atmospheric concentration of CO<sub>2</sub>.  
 152

153

## 154 2.2 The pCO<sub>2</sub>-Residual approach

155 We used the pCO<sub>2</sub>-Residual approach following Bennington et al. (2022a), which removes the  
 156 well-studied direct effect of temperature on pCO<sub>2</sub> from the LET model output before algorithm  
 157 processing. Temperature has both direct and indirect effects on surface ocean pCO<sub>2</sub>. The direct  
 158 effect of temperature, due to solubility and chemical equilibrium, is that an increase in temperature  
 159 directly causes an increase in pCO<sub>2</sub> (Takahashi et al., 1993). Indirectly, temperature changes are  
 160 associated with biological production and wintertime vertical mixing; and these processes tend to  
 161 result in opposing pCO<sub>2</sub> changes. To build reconstruction algorithms through the data-driven  
 162 training that occurs in ML, the statistics in all other algorithms developed to date must identify a  
 163 function that disentangles these competing effects of SST on pCO<sub>2</sub>. Here, the algorithm is assisted  
 164 by removing this known temperature effect, and it must therefore only learn the pCO<sub>2</sub> impacts

165 from biogeochemical drivers. The pCO<sub>2</sub>-Residual method leads to physically understandable  
166 connections between the input data and output (Bennington et al., 2022a), which mitigates to some  
167 degree ‘black box’ concerns typically associated with ML algorithms (Toms et al., 2020).  
168 Bennington et al. (2022a) demonstrate higher skill for reconstructions using pCO<sub>2</sub>-Residual as the  
169 target variable as opposed to pCO<sub>2</sub> (Figure S1 in Bennington et al., 2022a), indicating that the  
170 removal of the temperature-driven component enhances the performance of the method. Further,  
171 the pCO<sub>2</sub>-Residual method has been shown to perform slightly better against independent  
172 observations than other common mapping methods (Bennington et al., 2022a). A brief description  
173 is provided here, but for further details see Bennington et al. (2022a).

174 The temperature-driven component of pCO<sub>2</sub> (pCO<sub>2</sub>-T) is calculated using this equation:

$$175 \quad pCO_2-T = pCO_2^{mean} * \exp[0.0423 * (SST-SST^{mean})]$$

176 where pCO<sub>2</sub><sup>mean</sup> and SST<sup>mean</sup> is the long-term mean of surface ocean pCO<sub>2</sub> and temperature,  
177 respectively, using all 1°x1° grid cells from the testbed. Alternative sources of mean pCO<sub>2</sub> were  
178 assessed by Bennington et al. (2022a), but they found no significant impact on the test statistics or  
179 reconstructed pCO<sub>2</sub>. Once pCO<sub>2</sub>-T is determined, pCO<sub>2</sub>-Residual is calculated as the difference  
180 between pCO<sub>2</sub> and the calculated pCO<sub>2</sub>-T:

$$181 \quad pCO_2-Residual = pCO_2 - pCO_2-T$$

182 Prior to algorithm processing, pCO<sub>2</sub>-Residual values > 250 µatm and < -250 µatm from the  
183 testbed were filtered out targeting values that are not representative of the real ocean. The majority  
184 of the pCO<sub>2</sub>-Residual values that were filtered out correspond to high pCO<sub>2</sub>, above the maximum  
185 value in SOCAT (816 µatm; Stamell et al., 2020). The excluded data points (less than 0.2 % per  
186 member) mostly occurred in output from the CanESM2 model, and were restricted geographically,  
187 predominantly along the western coastline of South America.

188 The eXtreme Gradient Boosting method (XGB; Chen and Guestrin, 2016) is used to  
189 develop an algorithm that allows driver variables (i.e., SST, SSS, Chl-a, MLD, xCO<sub>2</sub>) to predict  
190 the pCO<sub>2</sub>-Residual (**Fig. 1**; Step 2). The pCO<sub>2</sub>-Residual and associated feature variables is split  
191 into validation, training and testing sets. The test and validation set each account for 20 % of the  
192 data, leaving 60 % for training. The validation set is used to optimize the algorithm

193 hyperparameters, which define the architecture of decision trees used in the model. The training  
194 set is used to build the decision trees in XGB, while the test set is used to evaluate the performance  
195 of the final algorithm. The XGB algorithm for this study used 4,000 decision trees with a maximum  
196 depth of 6 levels, and this was fixed for all experiments (see **Supplementary Text A**). For the  
197 final reconstruction of surface ocean pCO<sub>2</sub> across all space and time points, the previously  
198 calculated pCO<sub>2</sub>-T values are added back to the reconstructed pCO<sub>2</sub>-Residual (**Fig. 1**; Step 3).

199 The full XGB process, including 1) training/evaluating/testing and 2) reconstructing  
200 globally at a monthly resolution, was repeated individually for each LET member. This process  
201 provided therefore a total of 75 unique reconstruction vs. ‘model truth’ pairs, which can be  
202 statistically compared (**Fig. 1**; Step 4).

### 203 *2.3 Statistical Analysis in the Testbed*

204 The statistical comparisons between the test set and the reconstructions are equivalent to what  
205 would be derived using real-world data (‘seen’ values). Here, we calculate error statistics based on  
206 the full reconstruction (pCO<sub>2</sub> from all 1°x1° grid cells of the testbed, except for those masked or  
207 filtered out). In the full reconstruction, ~ 99 % of the data do not correspond to SOCAT or  
208 Saildrone USV observations used to train the algorithm (**Fig. S1**). Training data would ideally be  
209 removed before performance evaluation, but since the training data represent only ~ 1 %, the  
210 impact of not removing them is negligible (**Fig. S2**). A suite of statistical metrics can be used to  
211 compare the reconstruction to the ‘model truth’ in order to assess how well the algorithm can  
212 extrapolate from sparse data to full-field coverage (**Fig. 1**; Step 4). In this study, we focus on bias  
213 and root-mean-squared error (RMSE). Bias is calculated as ‘mean prediction – mean observation’  
214 (i.e., pCO<sub>2</sub> predicted by XGB subtracted by the pCO<sub>2</sub> ‘model truth’), and is a measure of over- or  
215 underestimation in the reconstructions. RMSE measures the magnitude of the predicted error and  
216 is calculated as the square root of the mean of the squared errors. We focus our discussion on the  
217 mean across 75 members of the testbed for bias and RMSE. The spread across testbed ensemble  
218 members is non-negligible and will be the focus of future work; here, we present the testbed spread  
219 primarily in the **Supplement**.

220

221



## 222 2.4 Overview of sampling patterns and model runs

223 First, we sampled target and driver variables from the LET based on sampling distributions  
224 equivalent to that of the SOCAT database ('SOCAT-baseline'). Then, we combined the 'SOCAT-  
225 baseline' with testbed output representing additional Saildrone USV coverage in the Southern  
226 Ocean. The additional Southern Ocean coverage was based on 1) the Sutton et al. (2021) sampling  
227 campaign from 2019 ('one-latitude' track) and 2) realistic potential future meridional USV  
228 observations ('zigzag' track) (see **Section 2.4.2; Fig. 2**). We performed a total of 10 experimental  
229 runs (**Table 1**). These represent different sampling approaches, including: 1) repeating USV  
230 sampling over a five- or ten-year period, 2) varying the number of USVs and thus the total number  
231 of monthly  $1^\circ \times 1^\circ$  observations, and 3) restricting all observations to southern hemisphere winter  
232 months. By comparing the different runs, we can assess whether or not certain targeted sampling  
233 strategies in the Southern Ocean can improve surface ocean  $p\text{CO}_2$  ML reconstructions. As  
234 discussed above, the LET runs to 2016 only (Gloege et al., 2021). Saildrone USV observations  
235 were therefore sampled from the testbed starting in year 2006 or 2007 (for the ten-year sampling)  
236 or 2012 (for the five-year sampling) until 2016, i.e., the final year of the testbed.

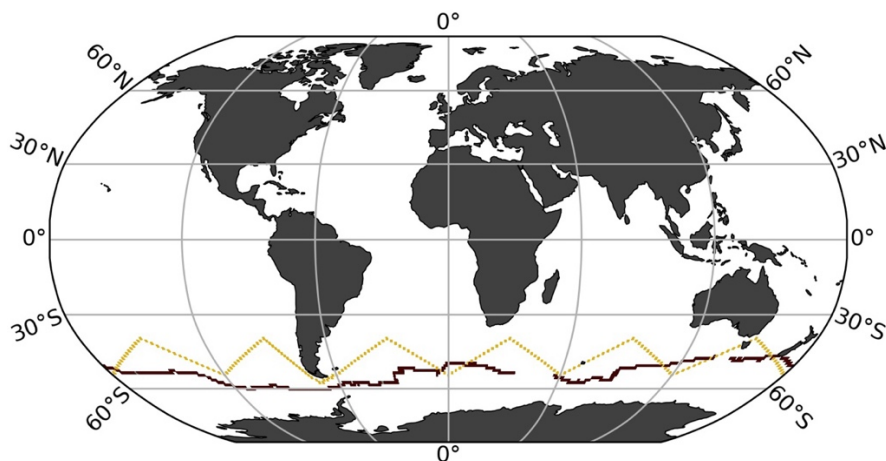
### 237 2.4.1 'One-latitude' runs

238 Six out of the ten experimental runs include the 'one-latitude' track (**Table 1**). The 2019 Saildrone  
239 USV journey (Sutton et al., 2021) covered an 8-month period, from January to August. Since the  
240 USV was recovered in early August, it did not cover the entire southern hemisphere winter (**Fig.**  
241 **S3**). We repeated this 'one-latitude' eight-month sampling pattern for five years ('5Y\_J-A'; 2,075  
242 observations) and ten years ('10Y\_J-A'; 4,150 observations). To evaluate year-round ('YR')  
243 coverage, the eight-month sampling period (January-August) was shifted by one month each year  
244 for ten years ('10Y\_YR'; 4,150 observations). To evaluate the impact of increased sampling, the  
245 2019 Saildrone USV track was repeated 12 times with incremental offsets of  $1^\circ$  from the original  
246 track, covering an additional  $6^\circ$  north and south (**Fig. S4**). This 'high-sampling'-run ('x13\_10Y\_J-  
247 A'; 44,250 observations) represents a total of 13 USVs. We also performed an additional 13 USV  
248 run, but including observations from southern hemisphere winter ('W') months only  
249 ('x13\_10Y\_W'; 25,395 observations). Finally, considering the cost of deploying 13 USVs, a  
250 downscaled 'multiple-USV-winter-only'-run was tested, including five USVs sampling over a

251 period of five years ('x5\_5Y\_W'; 5,022 observations). This run covers an additional 2° north and  
252 south from the original USV track.

#### 253 2.4.2 'Zigzag' runs

254 Four of the ten experimental runs represent realistic potential meridional sampling in the Southern  
255 Ocean ('zigzag' tracks; **Table 1**) as suggested by Djeutchouang et al. (2022). Saildrone USVs can  
256 operate at a speed capable of covering the spatial extent of meridional gradients in the Southern  
257 Ocean (Djeutchouang et al., 2022). However, Saildrone USVs are solar powered, and thus their  
258 range is restricted by the availability of solar radiation. To account for this and maintain a realistic  
259 sampling scenario, sampling occurs only to a maximum latitude of 55° S in these experiments.  
260 This alternative sampling pattern represents USVs sailing west to east in a north/south 'zigzag'  
261 pattern covering 40° S and 55° S for every 30° of longitude (**Fig. 2**). We created two scenarios.  
262 For the first scenario, every 30° of longitude from 40° S and 55° S is visited every three months  
263 within a single year as suggested by Lenton et al. (2006). Assuming an average Saildrone USV  
264 speed, this scenario represents four platforms equally spaced around the Southern Ocean. This  
265 sampling pattern was repeated for 10 years, with year-round coverage ('Zx4\_10Y\_YR'; 7,600  
266 observations), and for southern hemisphere winter months only ('Zx4\_10Y\_W'; 2,500  
267 observations). The second scenario represents a 'high-sampling' strategy, where every 30° of  
268 longitude from 40° S and 55° S is visited approximately monthly. This can be achieved by  
269 deploying 10 platforms equally spaced around the Southern Ocean running at an average Saildrone  
270 USV speed. This sampling pattern is repeated for five years, sampling year-round  
271 ('Z\_x10\_5Y\_YR'; 11,400 observations) and during southern hemisphere winter months only  
272 ('Z\_x10\_5Y\_W'; 3,800 observations).



273  
 274 **Figure 2:** Saildrone Uncrewed Surface Vehicle (USV) tracks representing the first circumnavigation around  
 275 Antarctica from 2019 in maroon ('one-latitude' track; Sutton et al., 2021) and an alternative virtual route with  
 276 meridional coverage ('zigzag' track).

Run name	SOCAT-baseline	5Y J-A	10Y J-A	10Y YR	x13 10Y J-A	x13 10Y W	x5 5Y W	Z x4 10Y YR	Z x4 10Y W	Z x10 5Y YR	Z x10 5Y W
<i>Saildrone track</i>	NA	One-lat	One-lat	One-lat	One-lat	One-lat	One-lat	Zigzag	Zigzag	Zigzag	Zigzag
<i>Years of sampling</i>	NA	5	10	10	10	10	5	10	10	5	5
<i>Duration of sampling</i>	NA	Jan-Aug	Jan-Aug	Year-round	Jan-Aug	SO winter	SO winter	Year-round	SO winter	Year-round	SO winter
<i>Additional observations</i>	NA	2,075	4,150	4,150	44,250	25,395	5,022	7,600	2,500	11,400	3,800
<i>Global coverage increase (%)</i>	NA	0.01	0.02	0.02	0.3	0.1	0.03	0.04	0.01	0.07	0.02
<b>Mean bias (µatm)</b>											
<i>Testbed period (1982-2016)</i>											
Globally	0.63	0.59	0.59	0.52	0.53	<b>0.39</b>	0.57	0.51	0.51	0.45	0.44
NORTH (35°N-90°N)	<b>0.11</b>	0.24	0.20	0.25	0.20	0.17	0.16	0.16	0.16	0.12	0.20
MID (35°S-35°N)	0.23	0.21	0.22	0.14	0.20	0.15	0.23	0.20	0.18	<b>0.13</b>	0.18
SOUTH (90°S-35°S)	1.4	1.3	1.2	1.1	1.1	<b>0.80</b>	1.2	1.1	1.1	1.0	0.87
SO winter months (JJA)	1.3	1.2	1.2	1.1	1.1	<b>0.90</b>	1.2	0.93	1.0	0.94	0.95
SO summer months (DJF)	0.070	0.11	0.15	0.10	0.15	<b>0.019</b>	0.11	0.25	0.073	0.16	0.066
<i>2006/2012-2016</i>											
Globally	0.51*	0.27	0.34	0.28	0.19	<b>0.03</b>	0.21	0.23	0.24	0.17	0.07
SOUTH (90°S-35°S)	1.6*	0.93	1.1	1.0	0.72	<b>0.37</b>	0.73	0.89	0.92	0.67	0.55
SOUTH (90°S-35°S) Jun, Jul, Aug	4.2*	2.6	2.7	2.8	2.2	1.8	2.5	1.8	2.4	<b>1.2</b>	2.0
<b>Mean RMSE (µatm)</b>											
<i>Testbed period (1982-2016)</i>											
Globally	11.8	11.7	11.8	11.7	11.7	11.6	11.7	<b>11.5</b>	11.6	<b>11.5</b>	11.6
NORTH (35°N-90°N)	<b>13.0</b>	<b>13.0</b>	<b>13.0</b>	<b>13.0</b>	<b>13.0</b>	<b>13.0</b>	13.1	<b>13.0</b>	<b>13.0</b>	<b>13.0</b>	<b>13.0</b>
MID (35°S-35°N)	11.7	11.7	11.7	11.7	11.7	11.7	11.7	11.7	11.7	11.7	11.7
SOUTH (90°S-35°S)	11.5	11.3	11.4	11.2	11.1	11.0	11.3	10.7	11.0	<b>10.6</b>	11.0
<i>2006/2012-2016</i>											
Globally	11.6*	11.6	11.4	11.3	11.3	11.2	11.6	<b>11.0</b>	11.2	11.1	11.4
SOUTH (90°S-35°S)	12.4*	11.1	11.0	10.7	10.6	10.4	10.9	10.0	10.6	<b>9.7</b>	10.6
SOUTH (90°S-35°S) Jun, Jul, Aug	12.0*	11.3	11.2	10.9	10.5	10.3	11.1	10.3	10.6	<b>9.6</b>	10.3

277  
 278 **Table 1.** Overview of the different sampling experiments tested in this study, and mean bias and RMSE (in µatm) for  
 279 various time periods, latitude bands for all runs. Bold values represent the best score for each category. 'One-lat' =  
 280 'one-latitude' track; incorporates the Saildrone USV route from Sutton et al. (2021). 'Zigzag' = potential meridional  
 281 sampling. 'Additional observations' = number of 1°x1° monthly Saildrone USV observations in addition to SOCAT.  
 282 J-A= January-August. YR = year-round. W = southern hemisphere winter. x4, x5, x10 and x13 = four, five, ten and  
 283 13 USVs. SO winter = Southern Ocean winter months, i.e., June, July, August and also including September. \*Average  
 284 value of the mean of 2006-2016 and 2012-2016. The global coverage increase was calculated based on the total  
 285 number of available 1982-2016 monthly 1°x1° observations from SOCAT (262,204 observations) and the Large  
 286 Ensemble Testbed (17,290,470 observations).

287  
 288 **2.5 Air-sea CO<sub>2</sub> flux**

289 To assess the global ocean carbon sink associated with our pCO<sub>2</sub> reconstructions, air-sea CO<sub>2</sub>  
 290 exchange was calculated for 1985 onward. Here, we computed air-sea CO<sub>2</sub> fluxes using the bulk

291 formulation with python package Seaflux.1.3.1 (<https://github.com/lukegre/SeaFlux>; Gregor et al.  
292 2021; Fay et al., 2021). We calculated global and Southern Ocean flux in the same manner for 1)  
293 the testbed ‘model truth’, 2) the ‘SOCAT-baseline’ and 3) the 10 experimental USV runs.

294 The net sea–air CO<sub>2</sub> flux was estimated using:

$$295 \text{ Flux} = k_w \cdot \text{sol} \cdot (\text{pCO}_2^{\text{ocn}} - \text{pCO}_2^{\text{atm}}) \cdot (1 - \text{ice})$$

296 where ‘ $k_w$ ’ is the gas transfer velocity, ‘sol’ is the solubility of CO<sub>2</sub> in seawater (in units of mol  
297 m<sup>-3</sup> μatm<sup>-1</sup>), ‘pCO<sub>2</sub><sup>ocn</sup>’ is the partial pressure of surface ocean carbon (in μatm), either from the  
298 ‘model truth’ or from the reconstructions, and pCO<sub>2</sub><sup>atm</sup> (in μatm) is the partial pressure of  
299 atmospheric CO<sub>2</sub> in the marine boundary layer. For GFDL, we used direct model output of  
300 pCO<sub>2</sub><sup>atm</sup>, while for CESM and CanESM2, pCO<sub>2</sub><sup>atm</sup> was calculated individually, as the product of  
301 surface xCO<sub>2</sub> and sea level pressure (the contribution of water vapor pressure was corrected for in  
302 CESM). Finally, to account for the seasonal ice cover in high latitudes, the fluxes were weighted  
303 by 1 minus the ice fraction (‘ice’), i.e., the open ocean fraction.

304 Winds have the largest impact on flux calculations (Fay et al., 2021), and temporally high-  
305 resolution output is not available for the LET. Monthly output is available, but this is not sufficient  
306 for the flux calculation due to the square dependency of wind speed (Wanninkhof, 2014). Given  
307 the necessity to use observed winds, for consistency, we use observations for all necessary  
308 variables for the flux calculation. Inputs to the calculation include EN4.2.2 salinity (Good et al.,  
309 2013), SST and ice fraction from NOAA Optimum Interpolation Sea Surface Temperature V2  
310 (OISSTv2) (Reynolds et al., 2002), and surface winds and associated wind scaling factor from the  
311 European Centre for Medium-Range Weather Forecasts (ECMWF ERA5 sea level pressure  
312 (Hersbach et al., 2020)). Results presented show the global and Southern Ocean (< 35° S) fluxes in  
313 units of Pg C yr<sup>-1</sup>.

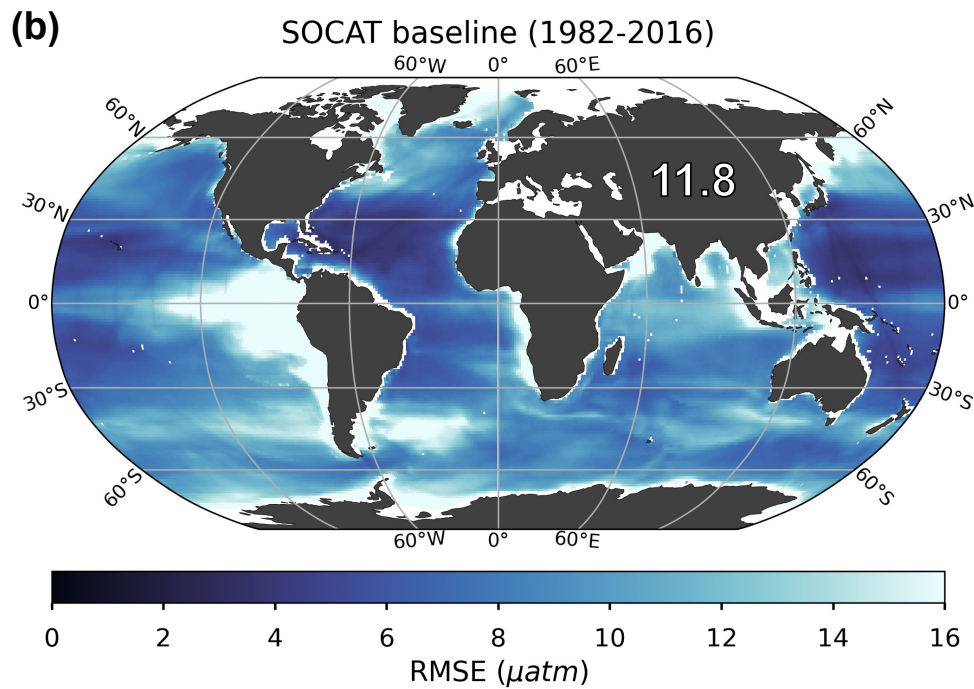
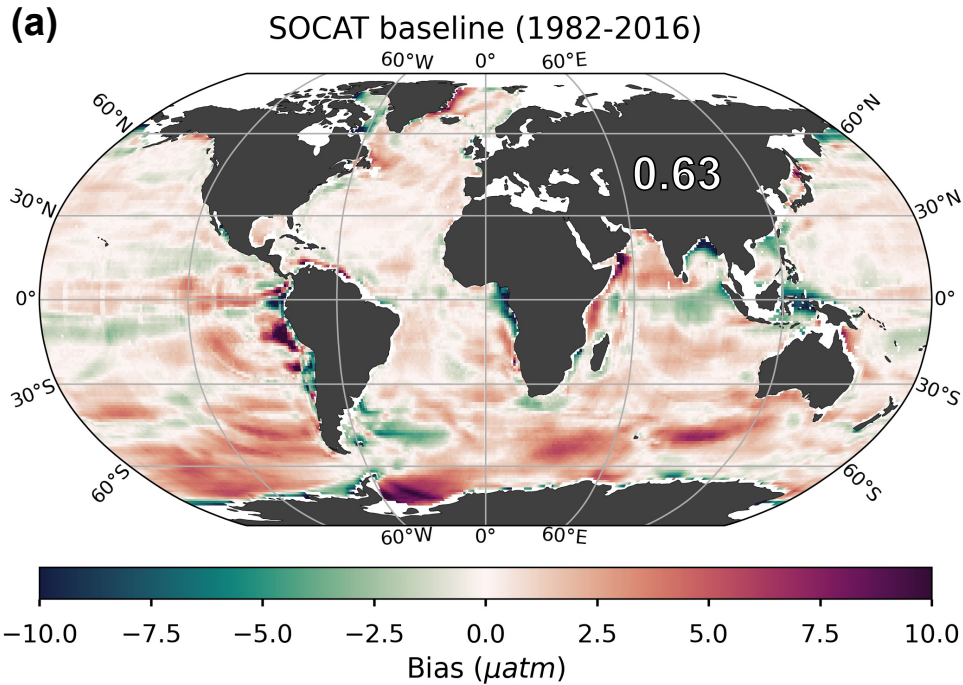
314 Note that, reconstructions of pCO<sub>2</sub> for the ‘SOCAT-baseline’ and the experimental USV  
315 runs are limited in their spatial extent to the open ocean (see **Sect. 2.1**; excluding coastal areas, the  
316 Arctic Ocean and marginal seas). The same mask was thus also applied when calculating the flux  
317 of the ‘model truth’, prior to comparison with the reconstructions.

318

### 319 3. Results

#### 320 3.1 Performance metrics for the ‘SOCAT-baseline’ reconstruction

321 The mean bias for the entire testbed period (i.e., 1982-2016) is 0.63  $\mu\text{atm}$  globally (**Fig. 3a**) and  
322 1.4  $\mu\text{atm}$  for the Southern Ocean ( $< 35^\circ \text{S}$ ; **Table 1**). Bias is much closer to zero for the mid-  
323 latitudes (between  $35^\circ \text{S}$  and  $35^\circ \text{N}$ ; 0.23  $\mu\text{atm}$ ) and northern latitudes ( $> 35^\circ \text{N}$ ; 0.11  $\mu\text{atm}$ ) (**Fig.**  
324 **3a**). There is a significant difference in bias considering southern hemisphere winter months (June,  
325 July, August) versus summer months (December, January, February), with a global mean bias (for  
326 1982-2016) of 1.3  $\mu\text{atm}$  compared to 0.07  $\mu\text{atm}$ , respectively (**Table 1**), due to the sparseness of  
327 SOCAT observations from the southern hemisphere during the harsh winter season (**Fig. S5a**).  
328 The mean RMSE for the entire testbed period (i.e., 1982-2016) is 11.8  $\mu\text{atm}$  globally (**Fig. 3b**) and  
329 11.5  $\mu\text{atm}$  for the Southern Ocean (**Table 1**). RMSE is highest in the Eastern Tropical and  
330 Southeastern Pacific Ocean and in the Southern Ocean, where the algorithm generally  
331 overestimates  $\text{pCO}_2$  (i.e., positive bias; **Fig. 3a**), with some exceptions in the Atlantic section. This  
332 is consistent with the areas significantly undersampled by SOCAT (**Fig. S5b**). Except for these  
333 areas, RMSE and bias is generally low (close to zero) in the open ocean, but show higher values  
334 along coastlines (**Fig. 3b**). The predicted  $\text{pCO}_2$  is thus more accurate in areas similar to and  
335 surrounding the SOCAT “observations” (i.e., monthly  $1^\circ \times 1^\circ$  grid cells equivalent to SOCAT  
336 coverage, but sampled from the LET). **Figure 3** shows mean bias and RMSE for the full  
337 reconstruction (see **Section 2.3**), but note that there is a statistically significant difference between  
338 the train and test set errors (**Fig. S6**). This indicates potential overfitting in our ML model (i.e.,  
339 higher errors for the ‘unseen’ reconstruction), and that further tuning of the hyperparameters could  
340 increase generalization skill (see **Supplementary Text A**).



341

342 **Figure 3:** Bias **(a)** and root-mean-squared error (RMSE) **(b)** for the ‘SOCAT-baseline’ (i.e., no USV) over the period  
 343 of 1982 through 2016. The global mean bias and RMSE is 0.63  $\mu\text{atm}$  and 11.8  $\mu\text{atm}$ , respectively. Note that only the  
 344 open ocean was considered in the reconstruction, so several areas were masked out prior to algorithm processing, such  
 345 as the Arctic Ocean, coastal areas and marginal seas (no data; white areas in figures).

346

347 *3.2 Reconstruction improvements with Saildrone USV additions*

348 Our presentation of global maps is limited to runs ‘x5\_5Y\_W’ (5,022 monthly 1°x1° observations)  
349 and ‘Z\_x4\_10Y\_YR’ (7,600 monthly 1°x1° observations). These runs were selected as they  
350 represent observational schemes that are realistic in the near-term future considering logistics and  
351 cost level, both non-meridional and meridional sampling, and different approaches to observing  
352 duration and seasonal coverage. For the remaining runs, equivalent maps can be found in the  
353 **Supplement**.

### 354 *3.2.1 Bias*

355 All Saildrone USV runs show a reduction in bias compared to the global mean 1982-2016  
356 ‘SOCAT-baseline’ (**Figs. 4a, S7**). The improvement in bias is mainly due to lower reconstructed  
357 pCO<sub>2</sub> values at southern latitudes, where the ‘SOCAT-baseline’ reconstruction generally  
358 overestimates pCO<sub>2</sub> (**Fig. 3a**). The global mean bias for ‘zigzag’ run ‘Z\_x4\_10Y\_YR’ is 0.51  
359 µatm, a higher improvement (19 %) over the ‘SOCAT-baseline’ compared to the ‘one-latitude’  
360 run ‘x5\_5Y\_W’ (11 % mean improvement; mean bias = 0.57 µatm;) (**Fig. 4a; Table 1**). Generally,  
361 the ‘zigzag’ runs show higher improvements from the ‘SOCAT-baseline’ (19-31 % improvement;  
362 resulting mean bias = 0.44-0.51 µatm) compared to the ‘one-latitude’ runs (7-19 % improvement;  
363 resulting mean bias = 0.52-0.59 µatm) (**Fig. S6; Table 1**). However, the ‘one-latitude’-run  
364 ‘x13\_10Y\_W’ that samples southern hemisphere winter months only, stands out with the lowest  
365 global mean (1982-2016) bias of 0.39 µatm, representing a 39 % mean improvement from the  
366 ‘SOCAT-baseline’ (**Table 1; Fig. S7**). This run, however, has three and five times more  
367 observations (25,395) than ‘Z\_x4\_10Y\_YR’ and ‘x5\_5Y\_W’, respectively.

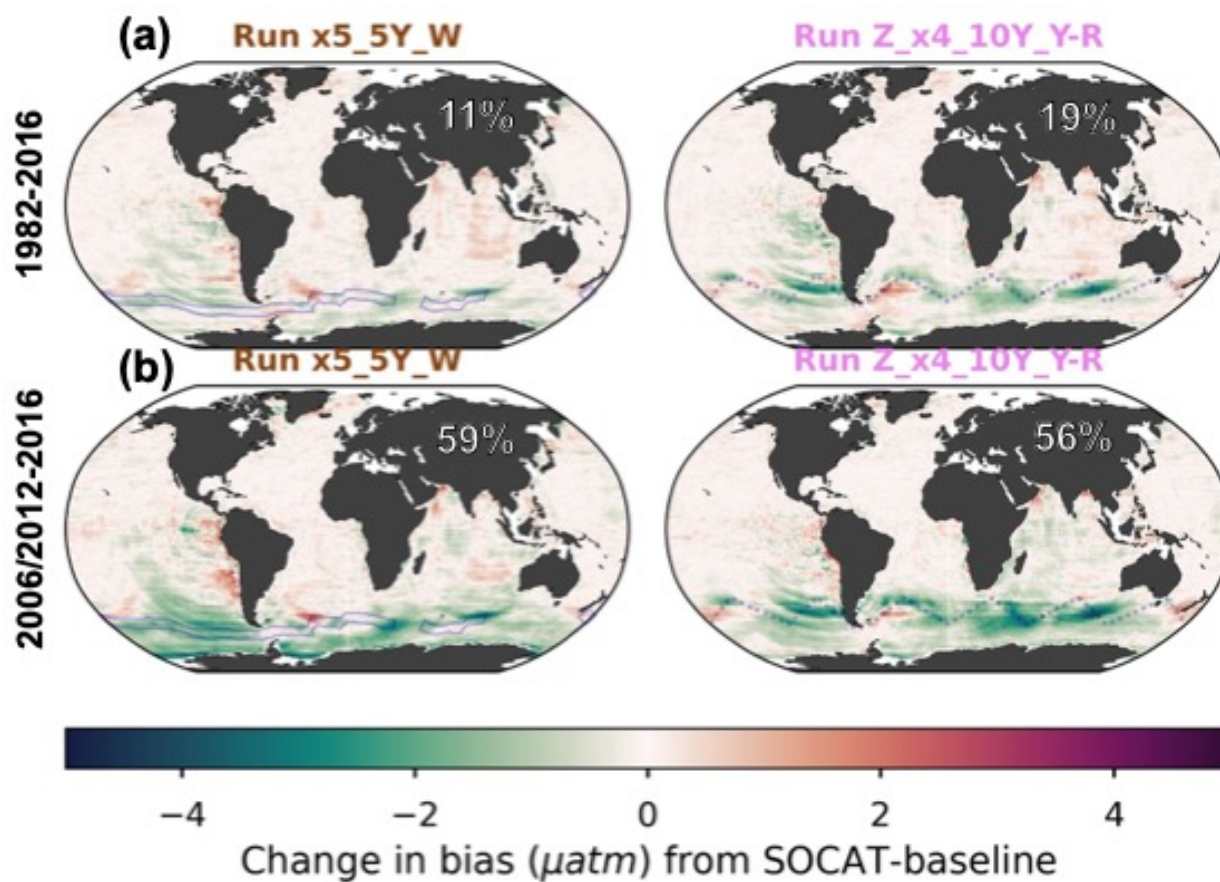
368 Compared to the entire testbed period, even larger improvements in global mean bias are  
369 shown for the period of Saildrone USV additions (2006-2016 and 2012-2016; **Figs. 4a vs. 4b,**  
370 **Figs. S7 vs. S8**). Compared to the ‘SOCAT-baseline’, run ‘x13\_10Y\_W’ results in a mean bias  
371 improvement of 95 %, while the remaining ‘one-latitude’ runs and the ‘zigzag’ runs show mean  
372 improvements up to 63 % and 85 %, respectively (**Fig. S8**). The spread in mean bias (2006/2012-  
373 2016) across the 75 testbed members for each experiment is shown in **Figure S9**.

374 Perhaps surprisingly, there is not a strong connection between the global or Southern Ocean  
375 mean bias and the number of added USV observations (**Fig. 5**). The ‘one-latitude’ ‘high-sampling’  
376 run ‘x13\_10Y\_J-A’ (44,250 observations) show similar mean bias or is outperformed by all

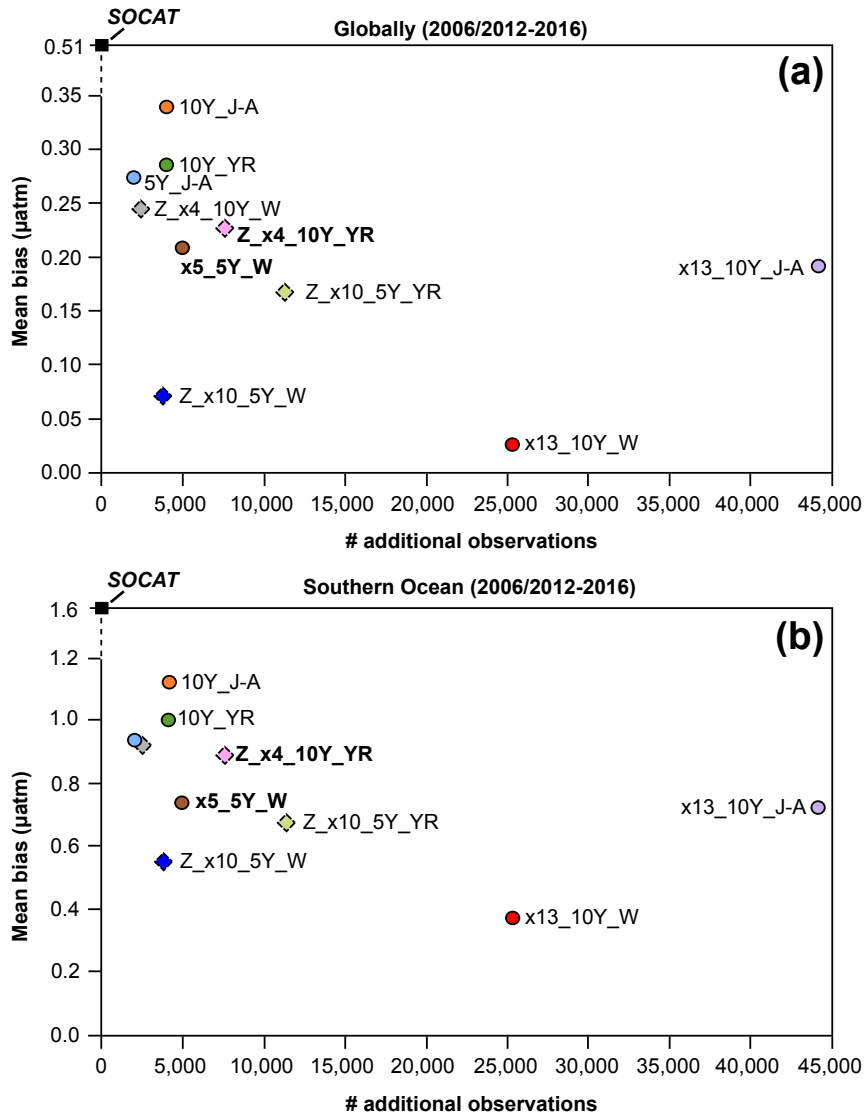
377 ‘zigzag’ runs as well as the ‘one-latitude’-runs that restrict sampling to southern hemisphere winter  
378 months (i.e., ‘x5\_5Y\_W’ and ‘x13\_10Y\_W’).

379           Considering the change in bias from year-to-year, the ‘SOCAT-baseline’ shows positive  
380 bias at all latitudes in the beginning of the testbed period, before improvement occurs around 1990  
381 (**Fig. 6a**). This is consistent with increasing SOCAT sampling with time for the period considered  
382 here (i.e., up to 2016; **Fig. S5c**). As SOCAT observations are biased towards the northern  
383 hemisphere (**Fig. S5a, b**), bias in the Southern Ocean ( $< 35^\circ$  S) increases significantly starting in  
384 the 2000s and remains high until the end of the testbed period (**Fig. 6a**). By adding USV sampling,  
385 bias in the Southern Ocean improves over the ‘SOCAT-baseline’ around year 2000 (**Fig. 6b-d**;  
386 **Fig. S10**), up to 6-12 years before to the introduction of additional samples in either 2006 or 2012.  
387 This improvement is shown for the majority of the 75 ensemble members (**Fig. S11**). Run  
388 ‘Z\_x10\_5Y\_W’, which has the lowest mean bias out of the ‘zigzag’ runs (**Fig. 5**), shows  
389 improvement even further back in time, until the beginning of the testbed period (**Fig. S10**). While  
390 the annual mean bias of the ‘zigzag’ runs varies rather consistently, there is a larger spread across  
391 the ‘one-latitude’ runs (**Fig. 6d**).

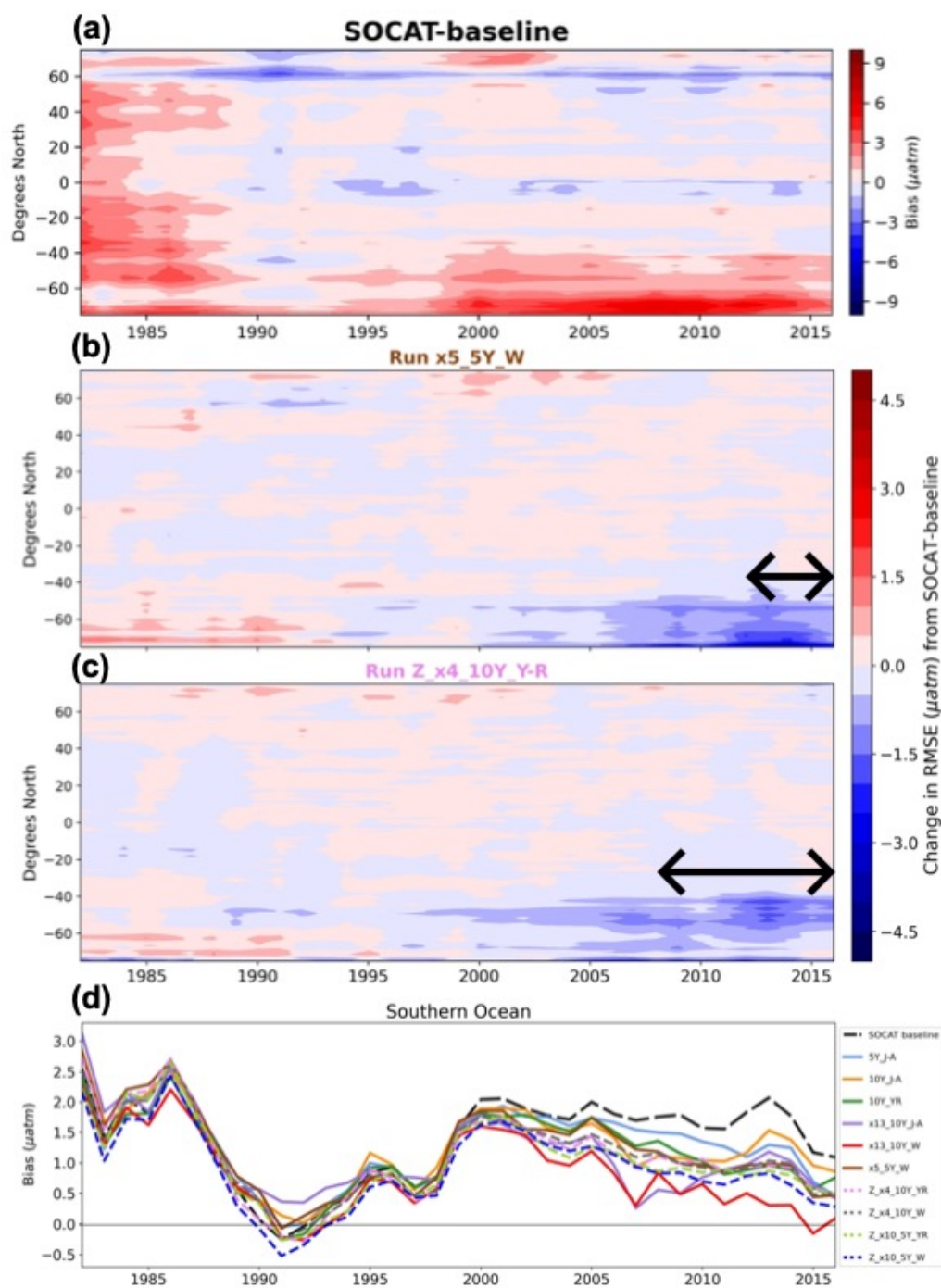




393  
 394 **Figure 4:** Change in bias when comparing run ‘x5\_5Y\_W’ and ‘Z\_x4\_10Y\_YR’ to the ‘SOCAT-baseline’  
 395 reconstruction, averaged over the duration of the testbed period (a; 1982-2016) and the period of USV additions (b;  
 396 2006-2012 or 2012-2016). The percent global improvement in absolute bias is shown on each panel. The USV  
 397 Sairdrone tracks are shown in blue.



398  
 399 **Figure 5:** Mean bias globally (a) and for the Southern Ocean (b) for the duration of Saildrone USV sampling (2006-  
 400 2016 or 2012-2016) for all runs presented in **Table 1**. Circles represent runs using the ‘one-latitude’ track, while  
 401 diamonds represent ‘zigzag’ runs. Runs highlighted in bold correspond to the two selected runs mapped in **Figure 4**,  
 402 **6**, **7** and **9**. Global (0.51 µatm) and Southern Ocean (1.6 µatm) bias values shown for the ‘SOCAT-baseline’ (black  
 403 squares) represent a mean of values for 2006-2016 (global = 0.52 µatm, S. Ocean = 1.63 µatm) and 2012-2016 (global  
 404 = 0.51 µatm, S. Ocean = 1.56 µatm). ‘# additional observations’ = number of monthly 1°x1° USV observations in  
 405 addition to SOCAT. Box plots illustrating the spread across the 75 ensemble members are shown in **Fig. S9**.



407

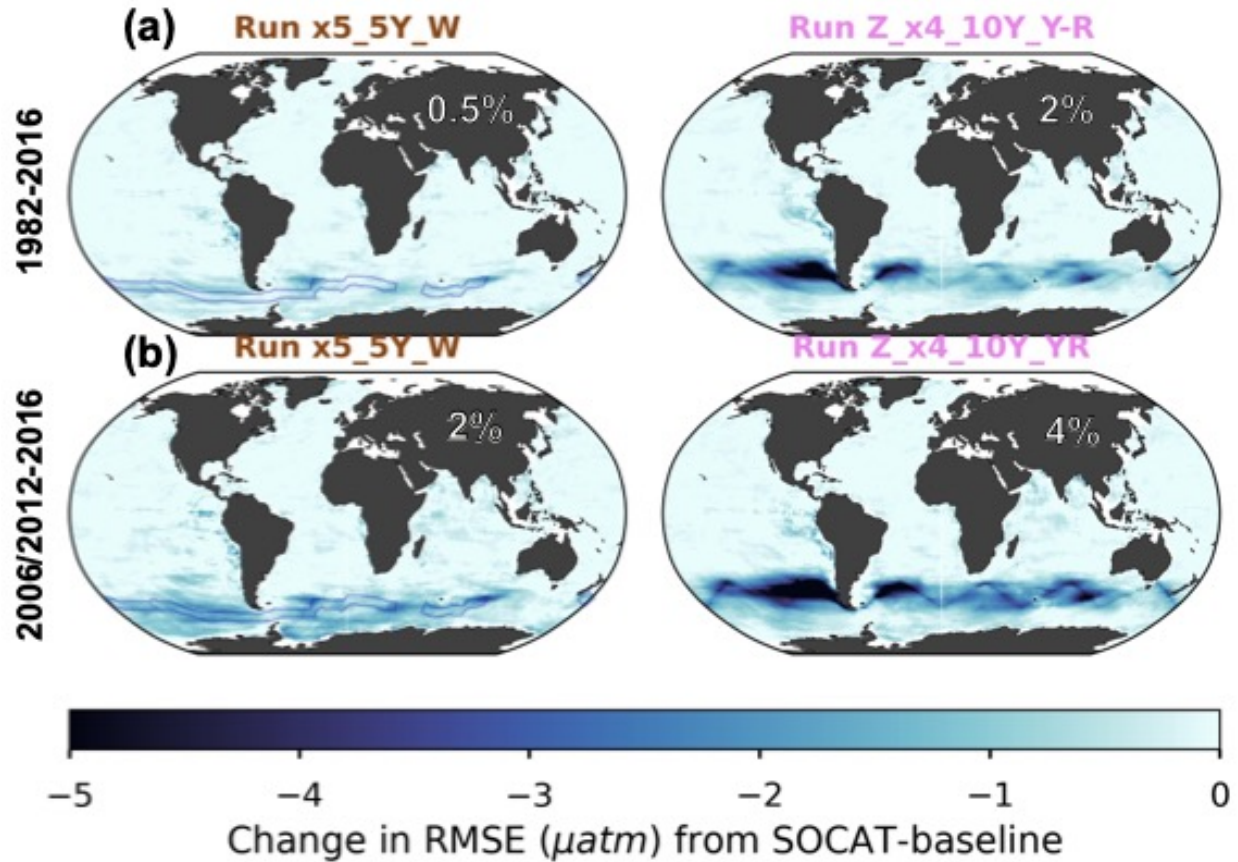
408 **Figure 6:** Zonal mean, annual mean Hovmöller of bias for the ‘SOCAT-baseline’ (a). Change in bias for run  
 409 ‘x5\_5Y\_W’ (b) and ‘Z\_x4\_10Y\_YR’ (c) compared to the ‘SOCAT-baseline’ shown in (a). Improvement in bias in  
 410 the Southern Ocean expands back in time well beyond the duration of USV additions for both runs (shown by arrows  
 411 on each panel). Annual mean bias for the Southern Ocean (> 35° S) for all runs (d).

412

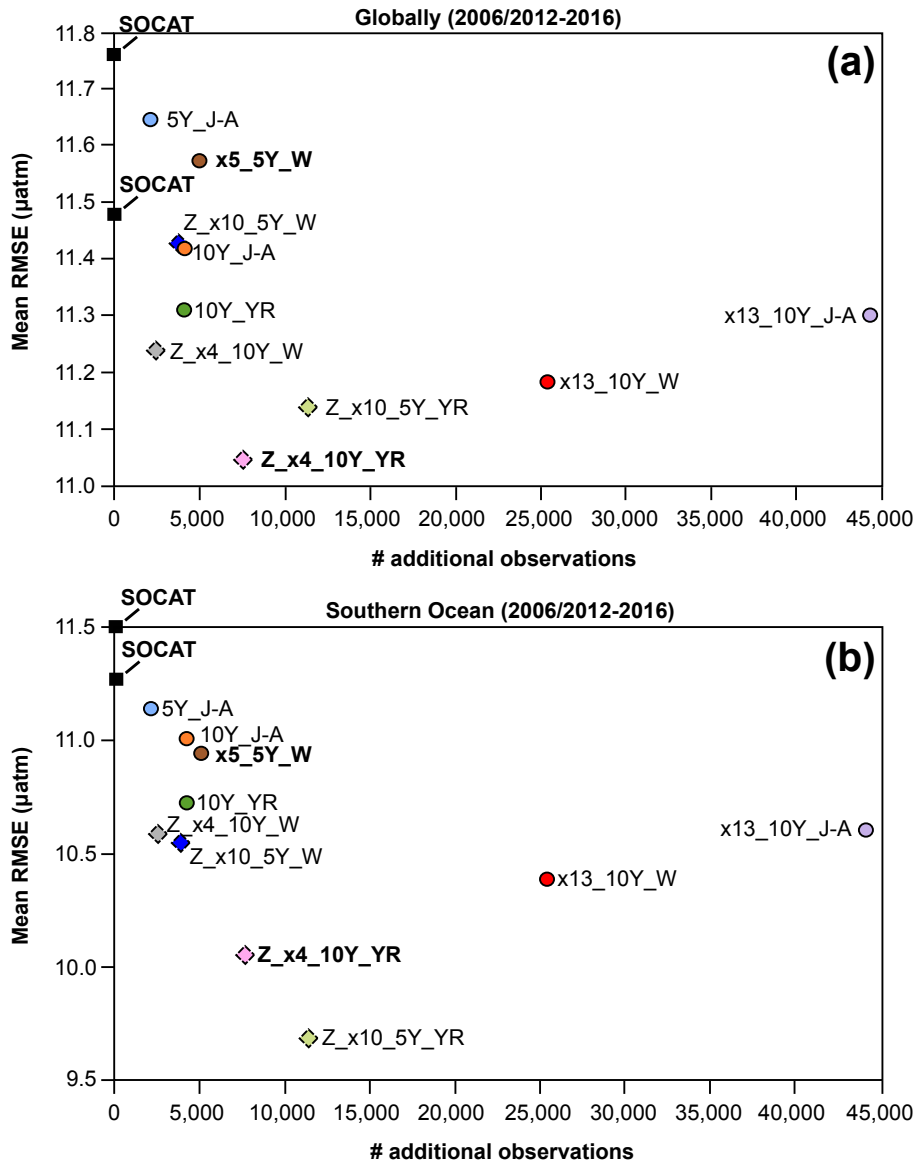
413 3.2.2 Root-mean squared error (RMSE)

414 Similar to bias, improvements in RMSE are most significant during the period of USV additions  
415 and within the Southern Ocean (**Fig. 7a** vs. **7b**). For the duration of USV additions, the ‘one-  
416 latitude’ runs show improvements in global mean RMSE of 1-3 % (0.1-1 % for 1982-2016), while  
417 the ‘zigzag’ runs show higher improvements between 2-5 % (1-3 % for 1982-2016) (**Figs. 7, S12,**  
418 **S13**). Mean RMSE is further reduced in the Southern Ocean by up to 16 %, and during southern  
419 hemisphere winter months (JJA) up to 21 % (run ‘Z\_x10\_5Y\_YR’; mean RMSE of 9.6  $\mu\text{atm}$ ;  
420 **Table 1**). There is minimal change in RMSE (or bias) during southern hemisphere summer months  
421 (DJF; **Fig. S14**). The two ‘zigzag’ runs sampling year-round (‘Z\_x4\_10Y\_YR’ and  
422 ‘Z\_x10\_5Y\_YR’) have the lowest RMSE values both globally and in the Southern Ocean (**Fig. 8**).  
423 The spread across the 75 testbed members for each experiment is shown in **Figure S15**.

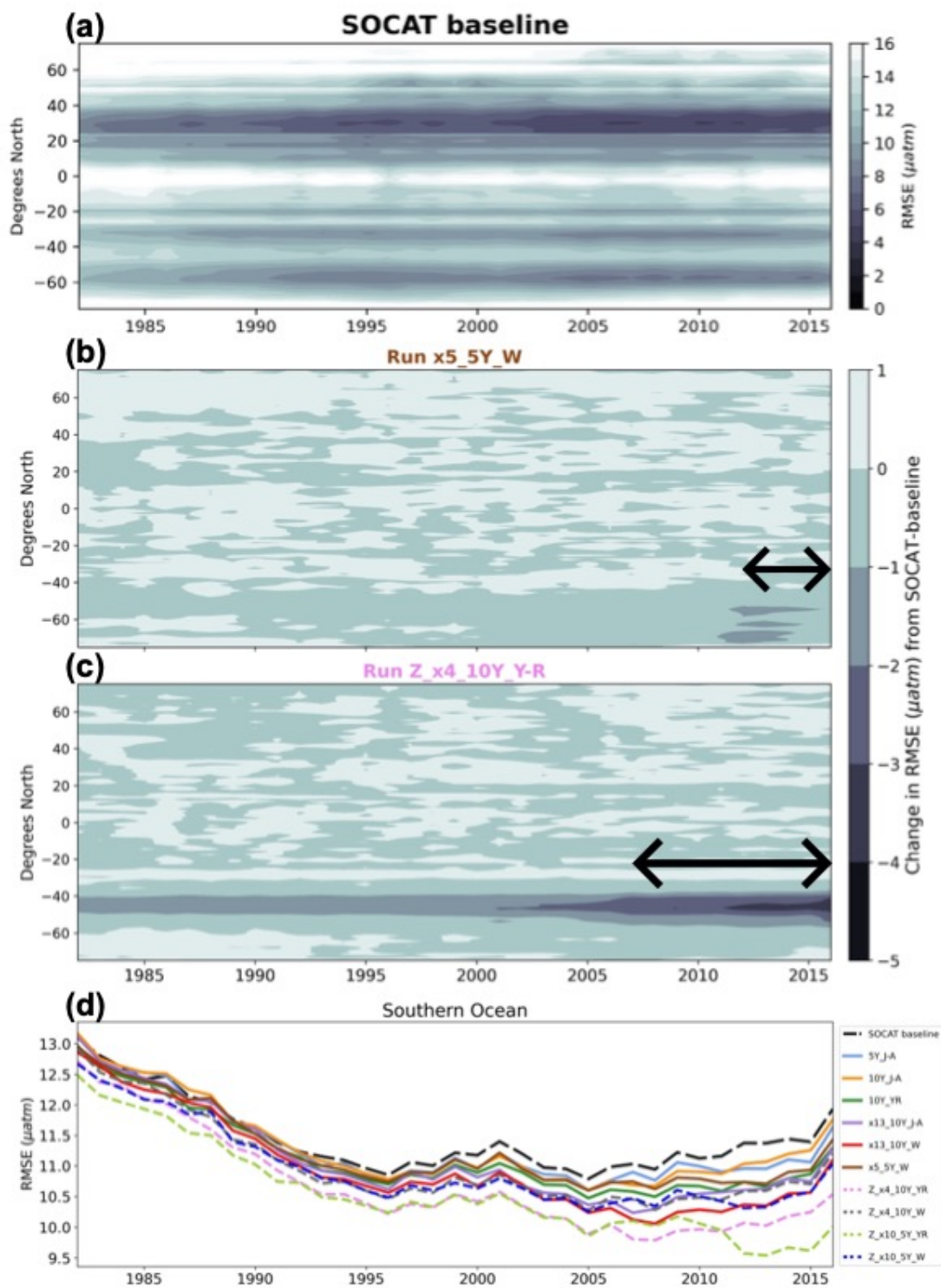
424 The ‘zigzag’ runs, as well as the ‘high-sampling’ ‘one-latitude’-runs (i.e., ‘x13\_10Y\_J-A’  
425 and ‘x13\_10Y\_W’), show improvements compared to the ‘SOCAT-baseline’ from the initiation  
426 of sampling (**Figs. 9, S16, S17**). The year-round ‘zigzag’ runs, however, show improvement in the  
427 Southern Ocean from the beginning of the testbed period (**Figs. 9c, d, S16**). RMSE improvements  
428 back in time are greater for all runs in the southern hemisphere winter months (**Fig. S18**).



429  
 430 **Figure 7:** Change in RMSE when comparing run 'x5\_5Y\_W' and 'Z\_x4\_10Y\_YR' to the 'SOCAT-baseline',  
 431 averaged over the duration of the testbed period (a; 1982-2016) and the period of Saildrone USV additions (b; 2006-  
 432 2012 or 2012-2016). The percent global improvement is shown on each panel.



433  
 434 **Fig. 8:** Mean RMSE globally (a) and for the Southern Ocean (< 35° S; b) for the duration of Saildrone USV sampling  
 435 (2006-2016 or 2012-2016) for all runs presented in **Table 1**. Circles represent runs using the ‘one-latitude’ track, while  
 436 diamonds represent ‘zigzag’ runs. Runs highlighted in bold correspond to the two selected runs mapped in **Figure 4,**  
 437 **6, 7** and **9**. RMSE values shown for the ‘SOCAT-baseline’ (black squares) represent a mean of values for 2006-2016  
 438 (global = 11.5 µatm, S. Ocean = 11.3 µatm) and 2012-2016 (global = 11.8 µatm, S. Ocean = 11.5 µatm). ‘# additional  
 439 observations’ = number of monthly 1°x1° USV observations in addition to SOCAT. Box plots illustrating the spread  
 440 across the 75 ensemble members are shown in **Fig. S15**.



443 **Figure 9:** Zonal mean, annual mean Hovmöller of RMSE for the ‘SOCAT-baseline’ (a). Change in RMSE for run  
 444 ‘x5\_5Y\_W’ (b) and ‘Z\_x4\_10Y\_YR’(c) compared to the ‘SOCAT-baseline’. Run ‘Z\_x4\_10Y\_YR’ shows  
 445 improvement in RMSE within the Southern Ocean, which expand well beyond the duration of Saildrone USV  
 446 additions (shown by arrow on panel). Annual mean RMSE for the Southern Ocean (> 35° S) for all runs (d).

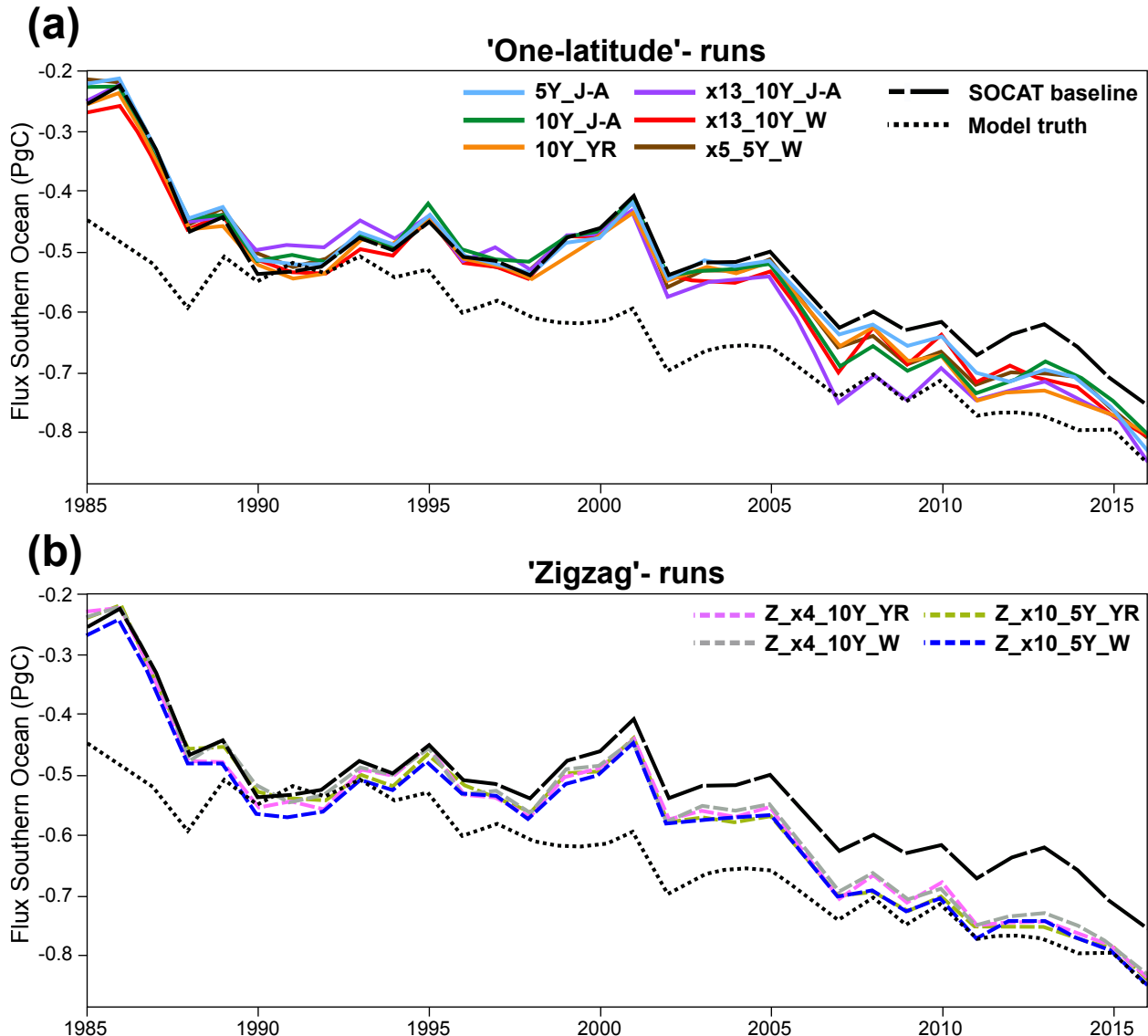
### 447 3.3 Impact on the air-sea CO<sub>2</sub> flux with Sairdrone USV additions

448 Air-sea flux was calculated in the same manner for both the ML reconstructions and the ‘model  
449 truth’, which allows for the isolation of the impact of different sampling strategies, as mediated by  
450 the pCO<sub>2</sub> reconstruction, on fluxes (see **Sect. 2.5**). These flux estimates are made to inform  
451 understanding of the errors that may exist in CO<sub>2</sub> flux estimates derived from pCO<sub>2</sub>  
452 reconstructions, and how new sampling could address these errors. Flux estimates represent the  
453 average of the 75 members of the LET in each case, and are not estimates of real-world fluxes.

454 Compared to the ‘model truth’, the ‘SOCAT-baseline’ reconstruction underestimates the  
455 global and Southern Ocean sink by 0.11-0.13 Pg C yr<sup>-1</sup> over 1982-2016 (**Fig. 10; Table S1**).  
456 Regardless of sampling pattern, adding Sairdrone USV observations increases both the global and  
457 Southern Ocean mean sink compared to the ‘SOCAT-baseline’ (**Figs. 10, S19**). The ‘one-latitude’  
458 runs show an increase of 0.01-0.03 Pg C yr<sup>-1</sup> (2-6 % strengthening) of the Southern Ocean sink  
459 (1982-2016), while the ‘zigzag’ runs lead to an even stronger sink by 0.04-0.06 Pg C yr<sup>-1</sup> (7-11 %  
460 strengthening) (**Table S2**). When averaging over the years of Sairdrone USV sampling addition  
461 (i.e., 2006-2012 and 2012-2016), the Southern Ocean sink increases up to 0.09 Pg C yr<sup>-1</sup> (14 %  
462 strengthening) for the ‘one-latitude’ runs and up to 0.1 Pg C yr<sup>-1</sup> (15 % strengthening) for the  
463 ‘zigzag’ runs (**Table S2**). These same features are found for the global ocean (**Fig. S19; Table**  
464 **S2**).

465 All of the ‘zigzag’ runs quite closely match both the global and Southern Ocean ‘model  
466 truth’ air-sea CO<sub>2</sub> flux for the duration of sample additions (**Figs. 10, S19**). Except for the first  
467 couple of years of sample addition for the ‘high-sampling’-run ‘x13\_10Y\_J-A’, none of the ‘one-  
468 latitude’ runs can match the ‘model truth’ air-sea CO<sub>2</sub> flux, instead they all underestimate the flux  
469 (**Figs. 10, S19**). The ‘zigzag’ runs have impact on the air-sea flux from an earlier date, starting to  
470 pull the results away from the ‘SOCAT-baseline’ and toward the ‘model truth’ already in the late-  
471 1990s, while the ‘one-latitude’ runs do the same about a decade later (**Figs. 10, S19**).





472 **Figure 10:** Southern Ocean (< 35° S) annually averaged air-sea CO<sub>2</sub> flux for the ‘SOCAT-baseline’ (black dashed  
 473 line), ‘model truth’ (black dotted line) ‘one-latitude’ runs (a; solid lines) and ‘zigzag’ runs (b; dashed lines).  
 474

475

476

477 **4. Discussion**

478 We have tested the pCO<sub>2</sub>-Residual reconstruction method with the Large Ensemble Testbed (LET)  
 479 to estimate its fidelity and understand how new samples could increase skill. We find that,  
 480 regardless of the chosen Saildrone USV sampling pattern, the reduction in mean bias and mean  
 481 RMSE compared to the ‘SOCAT-baseline’ is most prominent within the Southern Ocean (< 35°  
 482 S) during the period of which Saildrone USV observations were added (Figs. 4, 6, 7, 9). However,  
 483 it is important to mention that the additional Southern Ocean sampling also impacts (improves)

484 the pCO<sub>2</sub> reconstructions globally (**Figs. 5a, 8a**). Based on our experiments, a combination of  
485 factors improve global and Southern Ocean pCO<sub>2</sub> reconstructions, including the type of sampling  
486 pattern and seasonality of sampling, and to some extent, the number of additional observations.  
487 Importantly, increasing the number of observations or duration of sampling (5 vs. 10 years) is not  
488 the sole determining factor for improving the reconstructions (**Figs. 5, 8**). This is best demonstrated  
489 by the ‘high-sampling’-run ‘x13\_10Y\_J-A’ (44,250 observations), which does not provide  
490 significantly better reconstructions, or is even outperformed, by runs with 2-18 times fewer  
491 observations. The runs that produce lower mean RMSE do include data throughout southern  
492 hemisphere winter (**Fig. 8**). Run ‘x13\_10Y\_J-A’ does not include more than a few observations in  
493 the month of August, as it follows the temporal pattern of the real-world ‘one-latitude’ Saldron  
494 USV expedition (**Figs. S3, S4**; Sutton et al., 2021). The ‘one-latitude’ runs ‘10Y\_J-A’ and  
495 ‘10Y\_YR’ are directly comparable in terms of sample duration, spatial extent and number of  
496 observations (**Table 1**), but the latter, which covers all months, always shows lower mean RMSE  
497 and bias (**Figs. 5, 6d, 8, 9d**). These examples attest to the importance of addressing the issue of  
498 significant undersampling in the Southern Ocean during the winter season (**Fig. S5a**).

499 Another important comparison is the ‘one-latitude’-run ‘x5\_5Y\_W’ (5,022 observations)  
500 and ‘zigzag’-run ‘Z\_x10\_5Y\_W’ (3,800 observations) that both sample during southern  
501 hemisphere winter months over a five-year period (**Table 1**), where the ‘zigzag’-run consistently  
502 performs better even though it includes fewer observations (**Figs. 5, 8**). Most of the runs that  
503 perform similar to, or outperform, the above-mentioned ‘high-sampling’-run ‘x13\_10Y\_J-A’  
504 (44,250 observations), sample in a ‘zigzag’ pattern. Out of all 10 runs, the ‘year-round’ ‘zigzag’  
505 runs (‘Z\_x4\_10Y\_YR’ and ‘Z\_x10\_5Y\_YR’) are most able to reduce the mean error as shown by  
506 the lowest RMSE values (**Figs. 8, 9d**). A recent study performed similar sampling experiments as  
507 shown here, by comparing sampling from different types of autonomous platforms to a ‘SOCAT-  
508 baseline’ (Djeutchouang et al., 2022). They emphasized the importance of capturing the significant  
509 differences in pCO<sub>2</sub> that exist across meridional gradients during summer and winter months (up  
510 to 15  $\mu\text{atm}$ ; Djeutchouang et al., 2022). The meridional coverage provided by the ‘zigzag’ runs  
511 could explain why these runs generally outperform the ‘one-latitude’ runs in our study, and show  
512 significant reduction in both RMSE and bias, even though the global pCO<sub>2</sub> data density is raised  
513 by as little as 0.01-0.07 %.

514 The greatest reduction in mean bias out of all runs is shown by run ‘x13\_10Y\_W’ (**Figs.**  
515 **5, 6d**), which represents ‘one-latitude’ ‘high-sampling’ (i.e., 25,395 observations) during southern  
516 hemisphere winter months only. This sampling strategy seems thus to have a higher ability to  
517 reduce the ML model’s tendency to overestimate pCO<sub>2</sub> in the Southern Ocean compared to any of  
518 the meridional (‘zigzag’) runs. However, it should be noted that run ‘x13\_10Y\_W’ covers areas  
519 south of 55° S (**Fig. S4**), and its improvement in mean bias (and mean RMSE) is particularly  
520 prevalent at these high latitudes (e.g., **Figs. S8, S10, S13, S16**). Whether or not this run is, in fact,  
521 feasible with current or future technology is uncertain as parts of the southernmost tracks  
522 potentially cover the Southern Ocean ice zone (**Fig. S20**), and solar radiation for solar-powered  
523 platforms and sensors becomes very limited during winter south of 55° S. Furthermore, this  
524 particular sampling strategy requires 13 USVs, and so would be the most costly of the observing  
525 scenarios. Although run ‘x13\_10Y\_W’ demonstrates the highest reduction in mean bias out of all  
526 runs, the ‘zigzag’ runs still reduce absolute mean bias (for 2006/2012-2016) in the Southern Ocean  
527 by 44-65 % (vs. 77 % for run ‘x13\_10Y\_W’).

528 Overall, the ‘zigzag’ runs include significantly fewer observations, require fewer USVs,  
529 collect samples over the same duration, or even half the time as run ‘x13\_10Y\_W’, cover areas  
530 north of 55°S and within the ice-free zone, and show major improvement in the reconstruction of  
531 pCO<sub>2</sub>, attested to by reductions in both bias and RMSE. The ‘zigzag’ runs also closely match both  
532 the global and Southern Ocean ‘model truth’ air-sea CO<sub>2</sub> flux for the duration of sample additions  
533 (**Figs. 10, S19**). It also appears that the ‘zigzag’ runs generally have a greater impact on both the  
534 pCO<sub>2</sub> reconstruction and the air-sea flux further back in time, starting to deviate from the ‘SOCAT-  
535 baseline’ earlier compared to the ‘one-latitude’ runs (**Figs. 6, 9, 10, S10, S16, S18, S19**). Even the  
536 ‘zigzag’ scenarios with the least number of USVs (e.g., ‘Z\_x4\_10Y\_YR’) reduces Southern Ocean  
537 reconstruction absolute mean (2006-2016) bias and RMSE by up to 46 % and 11 %, respectively,  
538 and could provide a basis for realistic future Southern Ocean pCO<sub>2</sub> sampling campaigns.

539 The main motivation for improving surface ocean pCO<sub>2</sub> reconstructions is so that we can  
540 more accurately estimate the current and future oceanic uptake of anthropogenic carbon. The  
541 Southern Ocean is a significant carbon sink, but estimates of the air-sea CO<sub>2</sub> flux diverge  
542 substantially in this region (Takahashi et al., 2009; Landschützer et al., 2014, 2015; Rödenbeck et  
543 al., 2015; Williams et al., 2017; Gray et al., 2018; Gruber et al., 2019; Bushinsky et al., 2019; Long

544 et al., 2021; Fay and McKinley, 2021; Wu et al., 2022). Southern Ocean estimates incorporating  
545 observations from biogeochemical floats have shown a significantly weaker sink compared to  
546 those based only on observations from ships (Williams et al., 2017; Gray et al., 2018; Bushinsky  
547 et al., 2019). Bushinsky et al. (2019) and Hauck et al. (2023) performed similar sampling  
548 experiments as presented here, by comparing ML surface ocean pCO<sub>2</sub> reconstructions based on  
549 SOCAT vs. additional SOCCOM or ideal virtual floats. These studies showed that SOCAT  
550 sampling alone overestimates the CO<sub>2</sub> uptake in the Southern Ocean, and that additional floats  
551 reduce this overestimation, leading to a decreased (weakened) ocean carbon sink. In contrast, we  
552 find that the pCO<sub>2</sub>-Residual method underestimates the CO<sub>2</sub> uptake with only SOCAT sampling,  
553 and that adding USVs increased (strengthened) the Southern Ocean and global ocean sink by up  
554 to 0.1 Pg C yr<sup>-1</sup> (**Figs. 10, S19; Table S2**).

555         Going forward, additional studies are needed to better understand why these results suggest  
556 a different direction of the sink change with additional sampling. These differences could stem  
557 from the use of different reconstruction methods assessed. Hauck et al. (2023) used the MPI-SOM-  
558 FFN and CarboScope/Jena-MLS reconstruction methods, while we use the pCO<sub>2</sub>-Residual  
559 method. Another substantial difference between the studies is the models and numbers of ensemble  
560 members used as the testbed. Hauck et al. (2023) use a single hindcast model, while we use 25  
561 members each from three Earth System Models. We find substantial spread across these 75  
562 members (**Figs. S9 S15**), indicating that model structure and internal variability significantly  
563 impact results. Our study and Hauck et al. (2023) use different sampling masks and approaches  
564 for the calculation of fluxes, which could also be a factor. Targeted, coordinated studies using  
565 multiple reconstruction approaches with consistent testbed structures, sampling masks and  
566 experimental approaches are clearly needed (Rödenbeck et al., 2015). Despite this need for this  
567 additional work, studies do agree that additional Southern Ocean observations could significantly  
568 improve reconstructions of air-sea CO<sub>2</sub> fluxes.

569         What else can we learn using the model testbed? The ‘SOCAT-baseline’ demonstrates a  
570 weakening of the global and Southern Ocean carbon sink starting in the 1990s with a peak around  
571 year 2000 (**Figs. 10, S19**), which is in broad agreement with various data products using real-world  
572 SOCAT data (e.g., Gruber et al., 2019; Landschützer et al., 2015; Bushinsky et al., 2019;  
573 Bennington et al., 2022; Gloege et al., 2022). Peaks in bias and RMSE coincide in time with the

574 weakening sink (**Figs. 6d, 9d**). As shown by **Figure 10**, this ‘low sink’ is significantly exaggerated  
575 compared to the ‘model truth’. To better understand this discrepancy, we performed an additional  
576 experiment based on run ‘Z\_x10\_5Y\_YR’, but assumed sampling every year for the entire testbed  
577 period (i.e., 1982-2016). There is now a significant reduction in the temporal variability of  
578 reconstruction bias; with the additional 35-year USV sampling, the reconstructed Southern Ocean  
579 air-sea CO<sub>2</sub> flux closely matches the ‘model truth’ for the entire testbed duration (**Fig. S21**). This  
580 suggests that the large decadal variability of air-sea CO<sub>2</sub> fluxes since the 1980s, and the weak  
581 anomaly in the Southern Ocean carbon sink in the early 2000s (Le Quéré et al., 2007; Landschützer  
582 et al., 2015; Gruber et al., 2019; Bennington et al., 2022a,b; Friedlingstein et al., 2023), may be at  
583 least partially attributable to undersampling of the Southern Ocean. This is in agreement with the  
584 float sampling experiments performed by Hauck et al. (2023), attributing the strong decadal  
585 variability to sparse and skewed SOCAT data distributions. We will further explore this issue in  
586 future work. Still, this preliminary experiment suggests that interpretations of trends and variability  
587 of the global and Southern Ocean carbon sink should be considered with caution.

## 588 **5. Conclusions**

589 By using the Large Ensemble Testbed (LET), we show that targeted meridional and winter  
590 sampling in the Southern Ocean can improve global and Southern Ocean ML surface ocean pCO<sub>2</sub>  
591 reconstructions. Significant improvements are possible by raising the global pCO<sub>2</sub> data density by  
592 as little as 0.01-0.07 %. Further, we find that this modest amount of additional Saildrone USV  
593 sampling increases the global and Southern Ocean air-sea CO<sub>2</sub> flux by up to 0.1 Pg C yr<sup>-1</sup>, a  
594 quantity equivalent to 25 % of the uncertainty in the ocean carbon sink (0.4 Pg C yr<sup>-1</sup>;  
595 Friedlingstein et al., 2023). Our findings are consistent with previous studies suggesting that  
596 additional observations during southern hemisphere winter months and covering meridional  
597 gradients can reduce uncertainties and biases in the reconstructions (Lenton et al., 2006; Monteiro  
598 et al., 2010; Djeutchouang et al., 2022; Mackay et al., 2022). As opposed to other autonomous  
599 platform approaches, Saildrone USVs obtain in situ pCO<sub>2</sub> observations with uncertainties  
600 equivalent to the highest-quality observations collected by research ships ( $\pm 2 \mu\text{atm}$ ; Sabine et al.,  
601 2020; Sutton et al., 2021), and can operate at a high speed so that the spatial extent and seasonal  
602 cycle of meridional gradients can be covered. The approach of combining high-accuracy Saildrone  
603 USV and SOCAT observations represents thus a promising solution to improve future surface

604 ocean pCO<sub>2</sub> reconstructions and the accuracy of the ocean carbon sink. Lastly, we show that the  
605 large variability in bias, and the weakening of the global and Southern Ocean carbon sink in the  
606 2000s, may be partially an artefact of Southern Ocean undersampling.

#### 607 **Code availability**

608 Data analysis scripts will be made available in a GitHub repository upon publication.

#### 609 **Data availability**

610 The Large Ensemble Testbed is publicly available at  
611 [https://figshare.com/collections/Large\\_ensemble\\_pCO2\\_testbed/4568555](https://figshare.com/collections/Large_ensemble_pCO2_testbed/4568555).

612

#### 613 **Author contribution**

614 THH, GAM and AJS designed the experiments, and THH performed the simulations. THH, ARF  
615 and LG developed the code. THH and ARF calculated the air-sea fluxes. THH prepared the  
616 manuscript with contributions from all co-authors.

#### 617 **Competing interests**

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

#### 619 **Acknowledgements**

620 We acknowledge funding from NOAA through the Climate Observations and Monitoring Program  
621 (Award #NA20OAR4310340) and from NSF through the LEAP STC (Award #2019625). This is  
622 PMEL contribution 5549. We would also like to acknowledge and thank Val Bennington, Julius  
623 Busecke, Devan Samant and Abby Shaum for providing technical support, and Viviana Acquaviva  
624 for discussions regarding the manuscript. Lastly, we wish to thank two anonymous reviewers,  
625 whose contributions greatly improved the manuscript.

626

627

628 **References**

- 629
- 630 Bakker, D. C. E., Pfeil, B., Landa, C. S., Metzl, N., O'Brien, K. M., Olsen, A., Smith, K., Cosca,  
631 C., Harasawa, S., Jones, S. D., Nakaoka, S., Nojiri, Y., Schuster, U., Steinhoff, T., Sweeney, C.,  
632 Takahashi, T., Tilbrook, B., Wada, C., Wanninkhof, R., Alin, S. R., Balestrini, C. F., Barbero, L.,  
633 Bates, N. R., Bianchi, A. A., Bonou, F., Boutin, J., Bozec, Y., Burger, E. F., Cai, W.-J., Castle, R.  
634 D., Chen, L., Chierici, M., Currie, K., Evans, W., Featherstone, C., Feely, R. A., Fransson, A.,  
635 Goyet, C., Greenwood, N., Gregor, L., Hankin, S., Hardman-Mountford, N. J., Harlay, J., Hauck,  
636 J., Hoppema, M., Humphreys, M. P., Hunt, C. W., Huss, B., Ibáñez, J. S. P., Johannessen, T.,  
637 Keeling, R., Kitidis, V., Körtzinger, A., Kozyr, A., Krasakopoulou, E., Kuwata, A., Landschützer,  
638 P., Lauvset, S. K., Lefèvre, N., Lo Monaco, C., Manke, A., Mathis, J. T., Merlivat, L., Millero, F.  
639 J., Monteiro, P. M. S., Munro, D. R., Murata, A., Newberger, T., Omar, A. M., Ono, T., Paterson,  
640 K., Pearce, D., Pierrot, D., Robbins, L. L., Saito, S., Salisbury, J., Schlitzer, R., Schneider, B.,  
641 Schweitzer, R., Sieger, R., Skjelvan, I., Sullivan, K. F., Sutherland, S. C., Sutton, A. J., Tadokoro,  
642 K., Telszewski, M., Tuma, M., van Heuven, S. M. A. C., Vandemark, D., Ward, B., Watson, A.  
643 J., and Xu, S.: A multi-decade record of high-quality  $f\text{CO}_2$  data in version 3 of the Surface Ocean  
644  $\text{CO}_2$  Atlas (SOCAT), *Earth System Science Data*, 8, 383–413, [https://doi.org/10.5194/essd-8-383-](https://doi.org/10.5194/essd-8-383-2016)  
645 [2016](https://doi.org/10.5194/essd-8-383-2016), 2016.
- 646 Bakker, D. C. E., Alin, S. R., Becker, M., Bittig, H. C., Castaño-Primo, R., Feely, R. A., Gkritzalis,  
647 T., Kadono, K., Kozyr, A., Lauvset, S. K., Metzl, N., Munro, D. R., Nakaoka, S., Nojiri, Y., O'Brien,  
648 K. M., Olsen, A., Pfeil, Benjamin, P., Denis, S., Tobias, S., Kevin F., Sutton, A. J., Sweeney, C.,  
649 Tilbrook, B., Wada, C., Wanninkhof, R., Willstrand W. A., Akl, J., Apelthun, L. B., Bates, N.,  
650 Beatty, C. M., Burger, E. F., Cai, W., Cosca, C. E., Corredor, J. E., Cronin, M., Cross, J. N., De  
651 Carlo, E. H., DeGrandpre, M. D., Emerson, S. R., Enright, M. P., Enyo, K., Evans, W., Frangoulis,  
652 C., Fransson, A., García-Ibáñez, M. I., Gehrung, M., Giannoudi, L., Glockzin, M., Hales, B.,  
653 Howden, S. D., Hunt, C. W., Ibáñez, J. S. P., Jones, S. D., Kamb, L., Körtzinger, A., Landa, C.  
654 S., Landschützer, P., Lefèvre, N., Lo Monaco, C., Macovei, V. A., Maenner J. S., Meinig, C.,  
655 Millero, F. J., Monacci, N. M., Mordy, C., Morell, J. M., Murata, A., Musielewicz, S., Neill, .,  
656 Newberger, T., Nomura, D., Ohman, M., Ono, T., Passmore, A., Petersen, W., Petihakis, G.,  
657 Perivoliotis, L., Plueddemann, A. J., Rehder, G., Reynaud, T., Rodriguez, C., Ross, A. C.,  
658 Rutgeresson, A., Sabine, C. L., Salisbury, J. E., Schlitzer, R., Send, U., Skjelvan, I., Stamatakis, N.,

659 Sutherland, S. C., Sweeney, C., Tadokoro, K., Tanhua, T., Telszewski, M., Trull, T., Vandemark,  
660 D., van Ooijen, E., Voynova, Y. G., Wang, H., Weller, R. A., Whitehead, C., Wilson, D.: Surface  
661 Ocean CO<sub>2</sub> Atlas Database Version 2022 (SOCATv2022) (NCEI Accession 0253659), NOAA  
662 National Centers for Environmental Information [dataset], <https://doi.org/10.25921/1h9f-nb73>,  
663 2022.

664 Behncke, J., Landschützer, P. & Tanhua, T. A detectable change in the air-sea CO<sub>2</sub> flux estimate  
665 from sailboat measurements. *Scientific Reports*, 14, 3345, [https://doi.org/10.1038/s41598-024-](https://doi.org/10.1038/s41598-024-53159-0)  
666 [53159-0](https://doi.org/10.1038/s41598-024-53159-0), 2024.

667 Bennington, V., Galjanic, T., and McKinley, G. A.: Explicit Physical Knowledge in Machine  
668 Learning for Ocean Carbon Flux Reconstruction: The pCO<sub>2</sub>-Residual Method, *Journal of*  
669 *Advances in Modeling Earth Systems*, 14(10), <https://doi.org/10.1029/2021ms002960>, 2022a.

670 Bennington, V., Gloege, L., and McKinley, G. A.: Variability in the global ocean carbon sink from  
671 1959 to 2020 by correcting models with observations, *Geophysical Research Letters*, 49(14),  
672 <https://doi.org/10.1029/2022GL098632>, (2022b).

673 Bushinsky, S. M., Landschützer, P., Rödenbeck, C., Gray, A. R., Baker, D., Mazloff, M. R.,  
674 Resplandy, L., Johnson, K. S., and Sarmiento, J. L.: Reassessing Southern Ocean air-sea CO<sub>2</sub> flux  
675 estimates with the addition of biogeochemical float observations, *Global Biogeochemical Cycles*,  
676 33(11), 1370-1388, <https://doi.org/10.1029/2019GB006176>, 2019.

677 Chen, T., and Guestrin, C.: Xgboost: A scalable tree boosting system, In: *Proceedings of the 22nd*  
678 *ACM SIGKDD international conference on knowledge discovery and data mining* (pp. 785-794),  
679 <https://doi.org/10.1145/2939672.2939785>, 2016.

680 Denvil-Sommer, A., Gehlen, M., and Vrac, M.: Observation system simulation experiments in the  
681 Atlantic Ocean for enhanced surface ocean pCO<sub>2</sub> reconstructions, *Ocean Science*, 17, 1011-1030,  
682 <https://doi.org/10.5194/os-17-1011-2021>, 2021.

683 Deser, C., Phillips, A., Bourdette, V., and Teng, H.: Uncertainty in climate change projections: the  
684 role of internal variability, *Climate Dynamics*, 38, 527-546, [https://doi.org/10.1007/s00382-010-](https://doi.org/10.1007/s00382-010-0977-x)  
685 [0977-x](https://doi.org/10.1007/s00382-010-0977-x), 2012



686 Djeutchouang, L. M., Chang, N., Gregor, L., Vichi, M., and Monteiro, P. M. S.: The sensitivity of  
687  $p\text{CO}_2$  reconstructions to sampling scales across a Southern Ocean sub-domain: a semi-idealized  
688 ocean sampling simulation approach, *Biogeosciences*, 19, 4171-4195, [https://doi.org/10.5194/bg-](https://doi.org/10.5194/bg-19-4171-2022)  
689 [19-4171-2022](https://doi.org/10.5194/bg-19-4171-2022), 2022

690 Fay, A. R., Lovenduski, N. S., McKinley, G. A., Munro, D. R., Sweeney, C., Gray, A. R.,  
691 Landschützer, P., Stephens, B. B., Takahashi, T., and Williams, N.: Utilizing the Drake Passage  
692 Time-series to understand variability and change in subpolar Southern Ocean  $p\text{CO}_2$ ,  
693 *Biogeosciences*, 15(12), 3841-3855, <https://doi.org/10.5194/bg-15-3841-2018>, 2018.

694 Fay, A. R., and McKinley, G. A.: Observed regional fluxes to constrain modeled estimates of the  
695 ocean carbon sink, *Geophysical Research Letters*, 48(20), <https://doi.org/10.1029/2021GL095325>,  
696 2021.

697 Fay, A. R., Gregor, L., Landschützer, P., McKinley, G. A., Gruber, N., Gehlen, M., Iida, Y.,  
698 Laruelle, G. G., Rödenbeck, C., Roobaert, A., and Zeng, J.: SeaFlux: harmonization of air-sea  $\text{CO}_2$   
699 fluxes from surface  $p\text{CO}_2$  data products using a standardized approach, *Earth System Science Data*,  
700 13, 4693-4710, <https://doi.org/10.5194/essd-13-4693-2021>, 2021.

701 Friedlingstein, P., O'Sullivan, M., Jones, M. W., Andrew, R. M., Bakker, D. C. E., Hauck, J.,  
702 Landschützer, P., Le Quéré, C., Luijkx, I. T., Peters, G. P., Peters, W., Pongratz, J., Schwingshackl,  
703 C., Sitch, S., Canadell, J. G., Ciais, P., Jackson, R. B., Alin, S. R., Anthoni, P., Barbero, L., Bates,  
704 N. R., Becker, M., Bellouin, N., Decharme, B., Bopp, L., Brasika, I. B. M., Cadule, P.,  
705 Chamberlain, M. A., Chandra, N., Chau, T.-T.-T., Chevallier, F., Chini, L. P., Cronin, M., Dou,  
706 X., Enyo, K., Evans, W., Falk, S., Feely, R. A., Feng, L., Ford, D. J., Gasser, T., Ghattas, J.,  
707 Gkritzalis, T., Grassi, G., Gregor, L., Gruber, N., Gürses, Ö., Harris, I., Hefner, M., Heinke, J.,  
708 Houghton, R. A., Hurtt, G. C., Iida, Y., Ilyina, T., Jacobson, A. R., Jain, A., Jarníková, T., Jersild,  
709 A., Jiang, F., Jin, Z., Joos, F., Kato, E., Keeling, R. F., Kennedy, D., Klein Goldewijk, K., Knauer,  
710 J., Korsbakken, J. I., Körtzinger, A., Lan, X., Lefèvre, N., Li, H., Liu, J., Liu, Z., Ma, L., Marland,  
711 G., Mayot, N., McGuire, P. C., McKinley, G. A., Meyer, G., Morgan, E. J., Munro, D. R., Nakaoka,  
712 S.-I., Niwa, Y., O'Brien, K. M., Olsen, A., Omar, A. M., Ono, T., Paulsen, M., Pierrot, D., Pockock,  
713 K., Poulter, B., Powis, C. M., Rehder, G., Resplandy, L., Robertson, E., Rödenbeck, C., Rosan, T.  
714 M., Schwinger, J., Séférian, R., Smallman, T. L., Smith, S. M., Sospedra-Alfonso, R., Sun, Q.,

715 Sutton, A. J., Sweeney, C., Takao, S., Tans, P. P., Tian, H., Tilbrook, B., Tsujino, H., Tubiello, F.,  
716 van der Werf, G. R., van Ooijen, E., Wanninkhof, R., Watanabe, M., Wimart-Rousseau, C., Yang,  
717 D., Yang, X., Yuan, W., Yue, X., Zaehle, S., Zeng, J., and Zheng, B.: Global Carbon Budget 2023,  
718 Earth Syst. Sci. Data, 15, 5301–5369, <https://doi.org/10.5194/essd-15-5301-2023>, 2023.

719 Fyfe, J. C., Derksen, C., Mudryk, L., Flato, G. M., Santer, B. D., Swart, N. C., Molotch, N. P.,  
720 Zhang, X., Wan, H., Arora, V. K., Scinocca, J., and Jiao, Y.: Large near-term projected snowpack  
721 loss over the western United States, Nature communications, 8(1), 14996,  
722 <https://doi.org/10.1038/ncomms14996>, 2017.

723 Gloege, L., McKinley, G. A., Landschützer, P., Fay, A. R., Frolicher, T. L., and Fyfe, J. C.:  
724 Quantifying Errors in Observationally Based Estimates of Ocean Carbon Sink Variability, Global  
725 Biogeochemical Cycles, 35(4), <https://doi.org/10.1029/2020gb006788>, 2021.

726 Gloege, L., Yan, M., Zheng, T. and McKinley, G. A.: Improved quantification of ocean carbon  
727 uptake by using machine learning to merge global models and pCO<sub>2</sub> data, Journal of Advances in  
728 Modeling Earth Systems, 14(2), <https://doi.org/10.1029/2021MS002620>, 2022.

729

730 Good, S. A., Martin, M., and Rayner, N. A.: EN4: Quality controlled ocean temperature and  
731 salinity profiles and monthly objective analyses with uncertainty estimates, Journal of  
732 Geophysical Research Oceans, 118(12), 6704-6717, <https://doi.org/10.1002/2013JC009067>,  
733 2013.

734

735 Gray, A. R., Johnson, K. S., Bushinsky, S. M., Riser, S. C., Russell, J. L., Talley, L. D.,  
736 Wanninkhof, R., Williams, N. L., and Sarmiento, J. L.: Autonomous biogeochemical floats detect  
737 significant carbon dioxide outgassing in the high-latitude Southern Ocean, Geophysical Research  
738 Letters, 45(17), 9049-9057, <https://doi.org/10.1029/2018GL078013>, 2018.

739 Gregor, L., Lebehot, A. D., Kok, S., and Monteiro, P. M. S.: A comparative assessment of the  
740 uncertainties of global surface ocean CO<sub>2</sub> estimates using a machine-learning ensemble (CSIR-  
741 ML6 version 2019a) – have we hit the wall?, Geoscientific Model Development, 12, 5113-5136,  
742 <https://doi.org/10.5194/gmd-12-5113-2019>, 2019.

743 Gregor, L. and Fay, A. R.: Air-sea CO<sub>2</sub> fluxes for surface pCO<sub>2</sub> data products using a standardized  
744 approach, Zenodo [code], <https://doi.org/10.5281/zenodo.5482547>, 2021.

745 Gruber, N., Landschützer, P., and Lovenduski, N. S.: The variable Southern Ocean carbon sink,  
746 The Annual Review of Marine Science, 11, 159-86, [https://doi.org/10.1146/annurev-marine-](https://doi.org/10.1146/annurev-marine-121916-063407)  
747 [121916-063407](https://doi.org/10.1146/annurev-marine-121916-063407), 2019.

748 Hauck, J., Nissen, C., Landschützer, P., Rödenbeck, C., Bushinsky, S., and Olsen, A.: Sparse  
749 observations induce large biases in estimates of the global ocean CO<sub>2</sub> sink: and ocean model  
750 subsampling experiment, Philosophical Transactions Of the Royal Society A, 381:20220063,  
751 <https://doi.org/10.1098/rsta.2022.0063>, 2023.

752 Kay, J. E., Deser, C., Phillips, A., Mai, A., Hannay, C., Strand, G., Arblaster, J. M., Bates, S. C.,  
753 Danabasoglu, G., Edwards, J., Holland, M., Kushner, P., Lamarque, J-F., Lawrence, D., Lindsay,  
754 K., Middelton, A., Munoz, E., Neale, R., Oleson, K., Polvani, L., and Vertenstein, M.: The  
755 Community Earth System Model (CESM) large ensemble project: A community resource for  
756 studying climate change in the presence of internal climate variability, Bulletin of the American  
757 Meteorological Society, 96(8), 1333-1349, <https://doi.org/10.1175/BAMS-D-13-00255>, 2015.

758 Khatiwala, S., Primeau, F., and Hall, T.: Reconstruction of the history of anthropogenic CO<sub>2</sub>  
759 concentrations in the ocean, Nature, 462(7271), 346-349, <https://doi.org/10.1038/nature08526>,  
760 2009.

761 Landschützer, P., Gruber, N., Bakker, D. C. E., and Schuster, U.: Recent variability of the global  
762 ocean carbon sink, Global Biogeochemical Cycles, 28(9), 927-949,  
763 <https://doi.org/10.1002/2014GB004853>, 2014.

764 Landschützer, P., Gruber, N., Haumann, F. A., Rödenbeck, C., Bakker, D. C. E., Van Heuven, S.,  
765 Hoppema, M., Metzl, N., Sweeney, C., Takahashi, T., Brook, B., and Wanninkhof, R.: The  
766 reinvigoration of the Southern Ocean carbon sink, Science, 349(6253), 1221-1224,  
767 <https://doi.org/10.1126/science.aab2620>, 2015.

768 Landschützer, P., Tanhua, T., Behncke, J., and Keppler, L.: Sailing through the Southern Ocean  
769 seas of air-sea CO<sub>2</sub> flux uncertainty, Philosophical Transactions of the Royal Society A, 381,  
770 <https://doi.org/10.1098/rsta.2022.0064>, 2023.

771 Lenton, A. B., Matear, R. J., and Tilbrook, B.: Design of an observational strategy for quantifying  
772 the Southern Ocean uptake of CO<sub>2</sub>, *Global Biogeochemical Cycles*, 20, 1-11.  
773 <https://doi.org/10.1029/2005GB002620>, 2006.

774 Lenton, A. B., Tilbrook, B., Law, R. M., Bakker, D. C. E., Doney, S. C., Gruber, N., Ishii, M.,  
775 Hoppema, M., Lovenduski, N. S., Matear, R. J., McNeil, B. I., Metzl, N., Mikaloff Fletcher, S. E.,  
776 Monteiro, P. M. S., Rödenbeck, C., Sweeney, C., and Takahashi, T.: Sea-air CO<sub>2</sub> fluxes in the  
777 Southern Ocean for the period 1990-2009, *Biogeosciences*, 10, 4037-4054,  
778 <https://doi.org/10.5194/bg-10-4037-2013>, 2013.

779 Le Quéré, C., Rödenbeck, C., Buitenhuis, E. T., Conway, T. J., Lagenfelds, R., Gomez, A.,  
780 Labuschagne C., Ramonet, M., Nakazawa, T., Metzl, N., Gillett, N., and Heimann, M.: Saturation  
781 of the Southern Ocean CO<sub>2</sub> sink due to recent climate change, *Science*, 316(5832), 1735-1738,  
782 <https://doi.org/10.1126/science.1136188>, 2007.

783 Long, M. C., Stephens, B. B., McKain, K., Sweeney, C., Keeling, R. F., Kort, E. A., Morgan, E.  
784 J., Bent, J. D., Chandra, N., Chevallier, F., Commane, R., Daube, B. C., Krummel, P. B., Loh, Z.,  
785 Luijckx, I. T., Munro, D., Patra, P., Peters, W., Ramonet, M., Rödenbeck, C., Stavert, A., Tans, P.,  
786 and Wofsy, S. C.: Strong Southern Ocean carbon uptake evident in airborne observations, *Science*,  
787 374(6572), 1275-1280, <https://doi.org/10.1126/science.abi4355>, 2021.

788 Mackay, N., and Watson, A.: Winter air-sea CO<sub>2</sub> fluxes constructed from summer observations of  
789 the polar Southern Ocean suggest weak outgassing, *Journal of Geophysical Research: Oceans*,  
790 126(5), e2020JC016600, <https://doi.org/10.1029/2020JC016600>, 2021.

791 Mackay, N., Watson, A., Suntharalingam, P., Chen, Z., and Rödenbeck, C.: Improved winter data  
792 coverage of the Southern Ocean CO<sub>2</sub> sink from extrapolation of summertime observations,  
793 *Communications Earth & Environment*, 3, 265, <https://doi.org/10.1038/s43247-022-00592-6>,  
794 2022.

795 McKinley, G. A., Fay, A. R., Eddebbar, Y. A., Gloege, L., and Lovenduski, N. S.: External forcing  
796 explains recent decadal variability of the ocean carbon sink, *AGU Advances*, 1(2),  
797 e2019AV000149, <https://doi.org/10.1029/2019AV000149>, 2020.

798 Mongwe, N. P., Vichi, M., and Monteiro, P. M. S.: The seasonal cycle of  $p\text{CO}_2$  and  $\text{CO}_2$  fluxes in  
799 the Southern Ocean: diagnosing anomalies in CMIP5 Earth system models, *Biogeosciences*, 15(9),  
800 2851-2872, <https://doi.org/10.5194/bg-15-2851-2018>, 2018.

801 Monteiro, P. M. S., Gregor, L., Lévy, M., Maenner, S., Sabine, C. L., and Swart, S.: Intraseasonal  
802 variability linked to sampling alias in air-sea  $\text{CO}_2$  fluxes in the Southern Ocean, *Geophysical*  
803 *Research Letters*, 42(20), 8507-8514, <https://doi.org/10.1002/2015GL066009>, 2015.

804 Rodgers, K. B., Lin, J., and Frölicher, T. L.: Emergence of multiple ocean ecosystem drivers in a  
805 large ensemble suite with an Earth system model, *Biogeosciences*, 12(11), 3301-3320.  
806 <https://doi.org/10.5194/bg-12-3301-2015>, 2015.

807 Rödenbeck, C., Bakker, D. C. E., Gruber, N., Iida, Y., Jacobson, A. R., Jones, S., Landschützer,  
808 P., Metzl, N., Nakaoka, S., Olsen, A., Park, G.-H., Peylin, P., Rodgers, K. B., Sasse T. P., Schuster,  
809 U., Shutler, J. D., Valsala, V., Wannikkhof, R., and Zeng, J.: Data-based estimates of the ocean  
810 carbon sink variability – first results of the Surface Ocean  $p\text{CO}_2$  Mapping intercomparison  
811 (SOCOM), *Biogeosciences*, 12, 7251-7278, <https://doi.org/10.5194/bg-12-7251-2015>, 2015.

812 Sabine, C., Sutton, A., McCabe, K., Lawrence-Slavas, N., Alin, S, Feely, R., Jenkins, R., Maenner,  
813 S., Meinig, C., Thomas, J., van Ooijen, E., Passmore, A., and Tilbrook, B.: Evaluation of a new  
814 carbon dioxide system for autonomous surface vehicles, *Journal of Atmospheric and Oceanic*  
815 *Technology*, 37(8), 1305-1317, <https://doi.org/10.1175/JTECH-D-20-0010.1>, 2020.

816 Stamell, J., Rustagi, R. R., Gloege, L., and McKinley, G. A.: Strengths and weaknesses of three  
817 Machine Learning methods for  $p\text{CO}_2$  interpolation, *Geoscientific Model Development*  
818 *Discussions*[preprint], doi:10.5194/gmd-2020-311, 22 October 2020.

819 Sutton, A. J., Williams, N. L., and Tilbrook, B.: Constraining Southern Ocean  $\text{CO}_2$  flux uncertainty  
820 using uncrewed surface vehicle observations, *Geophysical Research Letters*, 48(3),  
821 e2020GL091748, <https://doi.org/10.1029/2020GL091748>, 2021.

822 Takahashi, T., Olafsson, J., Goddard, J. G., Chipman, D. W., and Sutherland, S. C.: Seasonal  
823 variation of  $\text{CO}_2$  and nutrients in the high-latitude surface oceans: A comparative study, *Global*  
824 *Biogeochemical Cycles*, 7(4), 843-878, <https://doi.org/10.1029/93GB02263>, 1993.

825 Takahashi, T., Sutherland, S. C., Wanninkhof, R., Sweeney, C., Feely, R. A., Chipman, D. W.,  
826 Hales, B., Friederich, G., Chavez, F., Sabine, C., Watson, A., Bakker, D. C. E., Schuster, U., Metzl,  
827 N., Yoshikawa-Inoue, H., Ishii, M., Midorikawa, T., Nojiri, Y., Körtzinger, A., Steinhoff, T.,  
828 Hoppema, M., Olafsson, J., Arnarson, T. S., Tilbrook, B., Johannessen, T., Olsen, A., Bellerby,  
829 R., Wong, C. S., Delille, B., Bates, N. R., and de Baar, H. J. W.: Climatological mean and decadal  
830 change in surface ocean pCO<sub>2</sub>, and net sea-air CO<sub>2</sub> flux over the global oceans, *Deep Sea Research*  
831 Part II: Topical Studies in Oceanography, 56(8-10), 554-557,  
832 <https://doi.org/10.1016/j.dsr2.2008.12.009>, 2009.

833 Toms, B. A., Barnes, E. A., and Ebert-Uphoff, I.: Physically interpretable neural networks for the  
834 geosciences: Applications to earth system variability, *Journal of Advances in Modeling Earth*  
835 *Systems*, 12(9), e2019MS002002, <https://doi.org/10.1029/2019MS002002>, 2020.

836 Wanninkhof, R.: Relationship between wind speed and gas exchange over the ocean revisited,  
837 *Limnology and Oceanography: Methods*, 12, 351-362, <https://doi.org/10.4319/lom.2014.12.351>,  
838 2014.

839 Williams, N. L., Juranek, L. W., Feely, R. A., Johnson, K. S., Sarmiento, J. L., Talley, L. D.,  
840 Dickson, A. G., Gray, A. R., Wanninkhof, R., Russell, J. L., Riser, S. C., and Takeshita, Y.:  
841 Calculating surface ocean pCO<sub>2</sub> from biogeochemical Argo floats equipped with pH: An  
842 uncertainty analysis, *Global Biogeochemical Cycles*, 31(3), 591-604,  
843 <https://doi.org/10.1002/2016GB005541>, 2017.

844 Wu, Y., Bakker, D. C. E., Achterberg, E. P., Silva, A. N., Pickup D. P., Li, X., Hartman, S.,  
845 Stappard, D., Qi, D., and Tyrrell, T.: Integrated analysis of carbon dioxide and oxygen  
846 concentrations as a quality control of ocean float data, *Communications Earth & Environment*, 3,  
847 92, <https://doi.org/10.1038/s43247-022-00421-w>, 2022.

848

849

850

851

852

853