

1 **Assessing improvements in global ocean pCO₂ machine learning reconstructions with**
2 **Southern Ocean autonomous sampling**

3 Thea H. Heimdal¹, Galen A. McKinley¹, Adrienne J. Sutton², Amanda R. Fay¹, Lucas Gloege³

4 ¹Columbia University and Lamont-Doherty Earth Observatory, Palisades, NY, USA

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

7 ³Open 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 CO₂. However, estimates
13 of the Southern Ocean air-sea CO₂ 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 pCO₂. Here, we use a Large Ensemble
16 Testbed (LET) of Earth System Models and the pCO₂-Residual reconstruction method to assess
17 improvements in pCO₂ reconstruction fidelity that could be achieved with additional autonomous
18 sampling in the Southern Ocean added to existing Surface Ocean CO₂ Atlas (SOCAT)
19 observations. The LET allows for a robust evaluation of the skill of pCO₂ reconstructions in space
20 and time through comparison to ‘model truth’. With only SOCAT sampling, Southern Ocean and
21 global pCO₂ 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 pCO₂. 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 CO₂ fluxes shown by SOCAT-only
28 sampling may be partially attributable to undersampling of the Southern Ocean.

29

Deleted: to

Deleted: ly

Deleted: e

Deleted: ,

Deleted: 6

Deleted: 9

Deleted: using

Deleted: ,

38 **1. Introduction**

39 The ocean plays an important role in mitigating climate change by sequestering anthropogenic
40 carbon emissions. From 1850 to 2023, the oceans have removed a total of 180 ± 35 Gt of carbon
41 (Friedlingstein et al., 2023). In order to fully understand the climate impacts from rising emissions,
42 it is essential to accurately quantify the air-sea CO₂ flux and the global ocean carbon sink in space
43 and time. The Surface Ocean CO₂ Atlas (SOCAT; Bakker et al., 2016) is the largest global
44 database of surface ocean CO₂ observations, with data starting in 1957. The main synthesis and
45 gridded products contain over 33 million high-quality direct shipboard measurements of fCO₂
46 (fugacity of CO₂) with an uncertainty of $< 5 \mu\text{atm}$ (Bakker et al., 2022). However, due to limited
47 resources for ocean observing, limited number of ships/routes, inaccessible regions and unsafe
48 waters, the database covers only about 1% of the global ocean at monthly $1^\circ \times 1^\circ$ spatial resolution
49 over the period of 1982-2023, and is highly biased towards the northern hemisphere.

50 Observation-based data products have been developed to estimate full-coverage surface
51 ocean pCO₂ across space and time by extrapolating to global coverage from these sparse SOCAT
52 observations (e.g., Landschützer et al., 2014; Rödenbeck et al., 2015; Gloege et al., 2022;
53 Bennington et al., 2022a,b). Most of these data products utilize machine learning (ML) algorithms
54 to estimate a non-linear function between a suite of driver variables (i.e., sea surface temperature
55 SST, sea surface salinity - SSS, mixed layer depth - MLD, Chlorophyll - Chl-a, xCO₂ -
56 atmospheric CO₂) and surface ocean pCO₂ (the target variable) where these are co-located. The
57 driver variables are proxies for processes influencing ocean pCO₂. Full-coverage driver variable
58 datasets are then processed through these ML algorithms to produce estimated global full-coverage
59 surface ocean pCO₂. Since the data products rely on pCO₂ observations to estimate functions
60 between the target and driver variables, data sparsity remains a fundamental limitation to this
61 technique.

62 It has been suggested that targeted sampling from autonomous platforms combined with
63 ships, filling in the state space of pCO₂, represents a path forward to improve surface ocean pCO₂
64 reconstructions (Bushinsky et al., 2019; Gregor et al., 2019; Gloege et al., 2021; Djeutchouang et
65 al., 2022; Landschützer et al., 2023; Hauck et al., 2023). One major obstacle, however, is that the
66 indirect pCO₂ estimates from floats have high uncertainties ($\pm 11.4 \mu\text{atm}$) and may be biased by
67 as much as $\sim 4 \mu\text{atm}$ (Bakker et al., 2016; Williams et al., 2017; Fay et al., 2018; Gray et al., 2018;

Deleted: against

Deleted: Since

Deleted: 7

Deleted: 2

Deleted: better constrain

Deleted: in

Deleted: T

Deleted: ;

Deleted: ;

Deleted: ;

Deleted: ;

Deleted: ;

Deleted: train the algorithms and thus produce these relationships...

Formatted: Subscript

Deleted: likely

83 Sutton et al., 2021; Mackay and Watson 2021; Wu et al 2022). ~~These large uncertainties and biases~~
84 ~~arise when pCO₂ is not measured directly as in the observations included in SOCAT, but is rather~~
85 ~~estimated using measurements of pH combined with a regression-derived alkalinity estimate~~
86 ~~(Williams et al., 2017; Gray et al., 2018). SOCAT includes only direct pCO₂ observations.~~ Biases
87 and uncertainties ~~may~~ have large impacts on global air-sea CO₂ flux estimates, given that the global
88 mean air-sea disequilibrium is only 5-8 μatm (McKinley et al., 2020). It is therefore critical that
89 bias and uncertainty corrections are well-constrained over different oceanic conditions and over
90 time.

91 Uncrewed Surface Vehicles (USVs), such as those manufactured and maintained by
92 Saildrone Inc., represent a new type of autonomous platform that can obtain direct pCO₂
93 observations with significantly lower uncertainties compared to other autonomous methods, and
94 equivalent to the highest-quality shipboard measurements contained in SOCAT (± 2 μatm; Sabine
95 et al., 2020; Sutton et al., 2021). Such improvements in sampling are critically important in the
96 undersampled Southern Ocean. This region is fundamental in terms of the ocean's ability to
97 remove carbon from the atmosphere, being responsible for ~ 40% of the global ocean uptake of
98 anthropogenic CO₂ (Khatiwala et al., 2009). Improved data coverage in the Southern Ocean
99 represents thus a major opportunity to advance our understanding of the global ocean carbon sink
100 (Lenton et al., 2006, 2013; Takahashi et al., 2009; Monteiro et al., 2015; Gregor et al., 2019; Gray
101 et al., 2018; Mongwe et al., 2018; Bushinsky et al., 2019; Sutton et al., 2021; Long et al., 2021;
102 Mackay et al., 2022; Wu et al., 2022; Landschützer et al., 2023; ~~Hauck et al., 2023~~). A combination
103 of SOCAT and Saildrone USV observations would include high-accuracy data from both the long
104 record and global coverage of ship tracks, and the expanded finer resolution of spatial and seasonal
105 coverage of the poorly sampled Southern Ocean. Importantly, Saildrone USVs are also able to
106 cover the spatial extent and seasonal cycle of the meridional gradients, which has been shown to
107 be critical in order to reduce errors in reconstructing surface ocean pCO₂ (Djeuthouang et al.,
108 2022). A combined approach, with autonomous samples such as those obtained from Saildrone
109 USVs, in addition to high-quality observations collected from ships, represents thus a promising
110 solution to improve surface ocean pCO₂ ML reconstructions.

111 Here, we assess to what extent surface ocean pCO₂ reconstructions can improve by
112 implementing the pCO₂-Residual machine learning (ML) reconstruction (Bennington et al., 2022a)

Formatted: Font: 12 pt

Formatted: Font: 12 pt

Formatted: Font: 12 pt

Formatted: Font: 12 pt

Formatted: Font: 12 pt

Formatted: Font: 12 pt

Formatted: Font: 12 pt, Subscript

Formatted: Font: 12 pt

Deleted: can

Deleted:

115 with the combined inputs of SOCAT and Saildrone USV coverage. However, instead of using **real-**
116 **world** observations, we sample the target (i.e., surface ocean pCO₂) and driver variables (i.e., SST,
117 SSS, MLD, Chl-a and xCO₂) from our Large Ensemble Testbed (LET) of Earth System Models
118 (ESMs) (e.g., Stamell et al., 2020; Gloege et al., 2021; Bennington et al., 2022a). There are two
119 major benefits of using a testbed compared to actual observations. First, in an ESM, **the** surface
120 ocean pCO₂ **field is provided precisely** at all **model** times and **1°x1° points**. Therefore, the pCO₂
121 reconstructed by the ML algorithm can be robustly evaluated in space and time against a known
122 ‘truth’ (i.e., ‘model truth’). The reconstruction evaluation is thus not limited to the availability of
123 sparse real-world ocean observations. Secondly, a testbed can be used to plan and evaluate the
124 impact of different sampling strategies on the reconstructed pCO₂. It is important to stress that, by
125 using a model testbed, we do not predict real-world surface ocean pCO₂ and air-sea CO₂ fluxes.
126 The goal here is to assess the accuracy with which an ML algorithm can reconstruct the ‘model
127 truth’ given inputs of samples consistent with real-world data coverage from the SOCAT database
128 and Saildrone USVs.

Deleted: actual

Deleted: known

Deleted: locations

129 By utilizing the observational coverage of SOCAT and Saildrone USV transects, we assess
130 to what extent the pCO₂-Residual method accurately reconstructs model surface ocean pCO₂ in
131 space and time. Additionally, we explore the timing, magnitude, duration and spatial extent of
132 Southern Ocean USV sample additions that most significantly improve the pCO₂ predictions.

133 2. Methods

134 2.1 The Large Ensemble Testbed (LET)

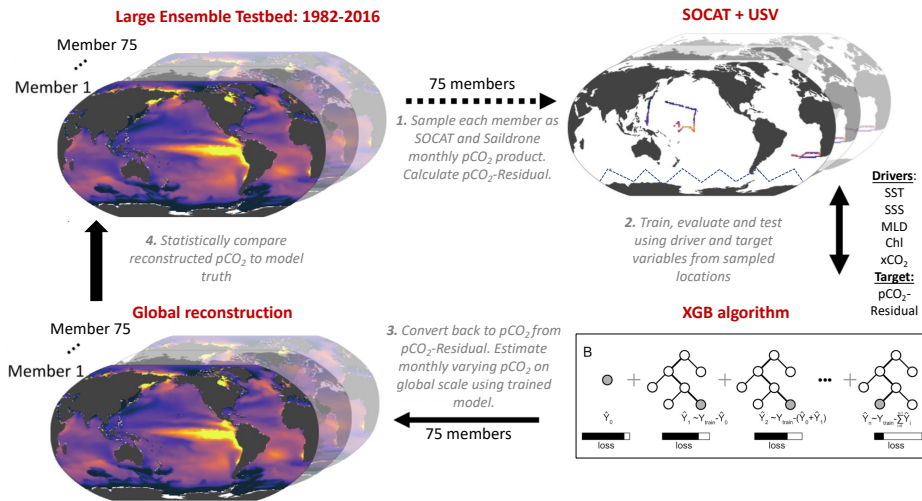
135 In this study, the Large Ensemble Testbed (LET) includes 25 members from three independent
136 initial-condition ensemble models (i.e., CanESM2, CESM-LENS and GFDL-ESM2M; Kay et al.,
137 2015; Rodgers et al., 2015; Fyfe et al., 2017), giving a total of 75 members within the testbed. We
138 do not use the MPI-GE model that was included in the past LET studies because its Southern
139 Ocean pCO₂ seasonality and decadal variability appear to be anomalously large (Gloege et al.,
140 2021; Fay and McKinley, 2021; Bennington et al., 2022a). Each individual Earth System Model
141 (ESM) is an imperfect representation of the actual Earth system, so the multiple Large Ensembles
142 are used to span different model structures and their representation of internal variability. Each
143 ensemble member undergoes the same external forcing (i.e., historical atmospheric CO₂ before

Deleted: s

148 2005 and Representative Concentration Pathway 8.5 through 2016, plus solar and volcanic
 149 forcing), but the spread across the ensemble members gives a unique trajectory of the ocean-
 150 atmosphere state over time, i.e., a different state of internal variability as well as the difference
 151 across models.

152 The LET used in this study includes monthly $1^\circ \times 1^\circ$ model output from 1982-2016 (Gloege
 153 et al., 2021). For each individual ensemble member of the LET, surface ocean $p\text{CO}_2$ and co-located
 154 driver variables (i.e., SST, SSS, Chl-a, MLD, $x\text{CO}_2$) were sampled monthly at a $1^\circ \times 1^\circ$ resolution,
 155 at times and locations equivalent to SOCAT and Saildrone USV observations (Fig. 1; Step 1).
 156 While the SOCAT observations were sampled from the testbed matching the actual years of
 157 sampling, the USV observations were sampled from the testbed starting in 2007 (for ten-year
 158 sampling) or 2012 (for five-year sampling) (see Sect. 2.4). As our focus is on reconstruction for
 159 the open ocean, testbed output for coastal areas, the Arctic Ocean ($>79^\circ\text{N}$) and marginal seas
 160 (Hudson Bay, Caspian Sea, Black Sea, Mediterranean Sea, Baltic Sea, Java Sea, Red Sea and Sea
 161 of Okhotsk) were removed prior to algorithm processing.

162



163 **Figure 1:** Schematic of the Large Ensemble Testbed (LET; modified from Gloege et al., 2021). **1:** Surface ocean
 164 $p\text{CO}_2$ from each of the 75 model members is sampled in space and time mimicking real-world SOCAT and Saildrone
 165 USV observations (see Fig. 2; Table 1; Section 2.5). Prior to algorithm processing, $p\text{CO}_2$ -Residual is calculated
 166 (Section 2.2). **2:** The $p\text{CO}_2$ -Residual (target variable) and co-located driver variables (i.e., SST, SSS, MLD, Chl,
 167

Deleted: year

Deleted: , i.e., the direct effect of temperature has been removed from the $p\text{CO}_2$ value

171 xCO₂) sampled from the testbed are processed by the XGBoost (XGB) algorithm (Section 2.3). 3: Based on the full-
172 coverage of driver variables, pCO₂-Residual is reconstructed globally. This process is repeated 75 times, individually
173 for every single testbed model member. The temperature component (pCO₂-T) is then added back to the pCO₂-
174 Residual for each value. 4: The globally reconstructed pCO₂ is evaluated against the 'model truth' at all 1°x1° grid
175 cells. SST = sea surface temperature. SSS = sea surface salinity. MLD = mixed layer depth. Chl = chlorophyll. xCO₂
176 = atmospheric concentration of CO₂.

177

178 2.2 The pCO₂-Residual approach

179 We used the pCO₂-Residual approach following Bennington et al. (2022a), which removes the
180 well-studied direct effect of temperature on pCO₂ from the LET model output before algorithm
181 processing. Temperature has both direct and indirect effects on surface ocean pCO₂. The direct
182 effect of temperature, due to solubility and chemical equilibrium, is that an increase in temperature
183 directly causes an increase in pCO₂ (Takahashi et al., 1993). Indirectly, temperature changes are
184 associated with biological production and wintertime vertical mixing; and these processes tend to
185 result in opposing pCO₂ changes. To build reconstruction algorithms through the data-driven
186 training that occurs in ML, the statistics in all other algorithms developed to date must identify a
187 function that disentangles these competing effects of SST on pCO₂. Here, the algorithm is assisted
188 by removing this known temperature effect, and it must therefore only learn the pCO₂ impacts
189 from biogeochemical drivers. The pCO₂-Residual method leads to physically understandable
190 connections between the input data and output (Bennington et al., 2022a), which mitigates to some
191 degree 'black box' concerns typically associated with ML algorithms (Toms et al., 2020). Further,
192 this method has been shown to perform better against independent observations than other
193 common observation-based products (Bennington et al., 2022a). A brief description is provided
194 here, but for further details see Bennington et al. (2022a).

195 The temperature-driven component of pCO₂ (pCO₂-T) is calculated using this equation:

$$196 \quad pCO_2-T = pCO_2^{mean} * \exp[0.0423 * (SST-SST^{mean})]$$

197 where pCO₂^{mean} and SST^{mean} is the long-term mean of surface ocean pCO₂ and temperature,
198 respectively, using all 1°x1° grid cells from the testbed. Once pCO₂-T is determined, pCO₂-
199 Residual is calculated as the difference between pCO₂ and the calculated pCO₂-T:

$$200 \quad pCO_2-Residual = pCO_2 - pCO_2-T$$

Deleted: Since we are using model testbed and not real-world observations, the

Deleted: can

Deleted: be

Deleted: , not just where observations are available

Deleted: ¶

Deleted: prior to

208 Prior to algorithm processing, pCO₂-Residual values > 250 μatm and < -250 μatm from the
209 testbed were filtered out, targeting values that are not representative of the real ocean. The majority
210 of the pCO₂-Residual values that were filtered out correspond to high pCO₂, above the maximum
211 value in SOCAT (816 μatm; Stamell et al., 2020). The excluded data points (less than 0.2 % per
212 member) mostly occurred in output from the CanESM2 model, and were restricted geographically,
213 predominantly along the western coastline of South America.

214 The eXtreme Gradient Boosting method (XGB; Chen and Guestrin, 2016) is used to
215 develop an algorithm that allows driver variables (i.e., SST, SSS, Chl-a, MLD, xCO₂) to predict
216 the pCO₂-Residual (Fig. 1; Step 2). The pCO₂-Residual and associated feature variables is split
217 into validation, training and testing sets. The test and validation set each account for 20 % of the
218 data, leaving 60 % for training. The validation set is used to optimize the algorithm
219 hyperparameters, which define the architecture of decision trees used in the model. The training
220 set is used to build the decision trees in XGB, while the test set is used to evaluate the performance
221 of the final algorithm. The XGB algorithm for this study used 4,000 decision trees with a maximum
222 depth of 6 levels, and this was fixed for all experiments. For the final reconstruction of surface
223 ocean pCO₂ across all space and time points, the previously calculated pCO₂-T values are added
224 back to the reconstructed pCO₂-Residual (Fig. 1; Step 3).

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

229 2.3 Statistical Analysis in the Testbed

230 The statistical comparisons between the test set and the reconstructions are equivalent to what
231 would be derived using real-world data (‘seen’ values). Here, we calculate error statistics based on
232 the full reconstruction (pCO₂ from all 1°x1° grid cells of the testbed, except for those masked or
233 filtered out). In the full reconstruction, ~99 % of the data do not correspond to SOCAT or
234 Saildrone USV observations used to train the algorithm (Fig. S1). Training data would ideally be
235 removed before performance evaluation, but since the training data represent only ~1 %, the
236 impact of not removing them is negligible (Fig. S2). A suite of statistical metrics can be used to

Deleted: to

Deleted: se

Deleted: generally

Formatted: Subscript

Deleted: Since we are using a testbed, we can also include comparisons on additional independent data, referred to as ‘unseen’ values, which represent the 1°x1° grid cells of the ensemble members that

Formatted: Not Highlight

Deleted:

Formatted: Font: Bold

Formatted: Font: Bold

245 compare the reconstruction to the ‘model truth’ in order to assess how well the algorithm can
246 extrapolate from sparse data to full-field coverage (**Fig. 1**; Step 4). In this study, we focus on bias
247 and root-mean-squared error (RMSE). Bias is calculated as ‘mean prediction – mean observation’
248 (i.e., pCO₂ predicted by XGB subtracted by the pCO₂ ‘model truth’), and is a measure of over- or
249 underestimation in the reconstructions. RMSE measures the magnitude of the predicted error and
250 is calculated as the square root of the mean of the squared errors. We focus our discussion on the
251 mean across 75 members of the testbed for bias and RMSE. The spread across testbed ensemble
252 members is non-negligible and will be the focus of future work; here, we present the testbed spread
253 primarily in the **Supplement**.

Formatted: Font: Bold

254 2.4 Overview of sampling patterns and model runs

255 First, we sampled target and driver variables from the LET based on sampling distributions
256 equivalent to that of the SOCAT database (‘SOCAT-baseline’). Then, we combined the ‘SOCAT-
257 baseline’ with testbed output representing additional Sairdrone USV coverage in the Southern
258 Ocean. The additional Southern Ocean coverage was based on 1) the Sutton et al. (2021) sampling
259 campaign from 2019 (‘one-latitude’ track) and 2) realistic potential future meridional USV
260 observations (‘zigzag’ track) (see Section 2.4.2; Fig. 2). We performed a total of 10 experimental
261 runs (**Table 1**). These represent different sampling approaches, including: 1) repeating USV
262 sampling over a five- or ten-year period, 2) varying the number of USVs and thus the total number
263 of monthly 1°x1° observations, and 3) restricting all observations to southern hemisphere winter
264 months. By comparing the different runs, we can assess whether or not certain targeted sampling
265 strategies in the Southern Ocean can improve surface ocean pCO₂ ML reconstructions. As
266 discussed above, the LET runs to 2016 only (Gloege et al., 2021). Sairdrone USV observations
267 were therefore sampled from the testbed starting in year 2006 or 2007 (for the ten-year sampling)
268 or 2012 (for the five-year sampling) until 2016, i.e., the final year of the testbed.

Deleted:

Deleted:

Formatted: Font: Bold

269 2.4.1 ‘One-latitude’ runs

270 Six out of the ten experimental runs include the ‘one-latitude’ track (**Table 1**). The 2019 Sairdrone
271 USV journey (Sutton et al., 2021) covered an 8-month period, from January to August. Since the
272 USV was recovered in early August, it did not cover the entire southern hemisphere winter (**Fig.**
273 **S3**). We repeated this ‘one-latitude’ eight-month sampling pattern for five years (‘5Y_J-A’; 2,075

Deleted: 1

277 observations) and ten years ('10Y_J-A'; 4,150 observations). To evaluate year-round ('YR')
278 coverage, the eight-month sampling period (January-August) was shifted by one month each year
279 for ten years ('10Y_YR'; 4,150 observations). To evaluate the impact of increased sampling, the
280 2019 Sairdron USV track was repeated 12 times with incremental offsets of 1° from the original
281 track, covering an additional 6° north and south (Fig. S4). This 'high-sampling'-run ('x13_10Y_J-
282 A'; 44,250 observations) represents a total of 13 USVs. We also performed an additional 13 USV
283 run, but including observations from southern hemisphere winter ('W') months only
284 ('x13_10Y_W'; 25,395 observations). Finally, considering the cost of deploying 13 USVs, a
285 downscaled 'multiple-USV-winter-only'-run was tested, including five USVs sampling over a
286 period of five years ('x5_5Y_W'; 5,022 observations). This run covers an additional 2° north and
287 south from the original USV track.

288 2.4.2 'Zigzag' runs

289 Four of the ten experimental runs represent realistic potential meridional sampling in the Southern
290 Ocean ('zigzag' tracks; Table 1) as suggested by Djeutchouang et al. (2022). Sairdron USVs can
291 operate at a speed capable of covering the spatial extent of meridional gradients in the Southern
292 Ocean (Djeutchouang et al., 2022). However, Sairdron USVs are solar powered, and thus their
293 range is restricted by the availability of solar radiation. To account for this and maintain a realistic
294 sampling scenario, sampling occurs only to a maximum latitude of 55° S in these experiments.
295 This alternative sampling pattern represents USVs sailing west to east in a north/south 'zigzag'
296 pattern covering 40° S and 55° S for every 30° of longitude (Fig. 2). We created two scenarios.
297 For the first scenario, every 30° of longitude from 40° S and 55° S is visited every three months
298 within a single year as suggested by Lenton et al. (2006). Assuming an average Sairdron USV
299 speed, this scenario represents four platforms equally spaced around the Southern Ocean. This
300 sampling pattern was repeated for 10 years, with year-round coverage ('Zx4_10Y_YR'; 7,600
301 observations), and for southern hemisphere winter months only ('Zx4_10Y_W'; 2,500
302 observations). The second scenario represents a 'high-sampling' strategy, where every 30° of
303 longitude from 40° S and 55° S is visited approximately monthly. This can be achieved by
304 deploying 10 platforms equally spaced around the Southern Ocean running at an average Sairdron
305 USV speed. This sampling pattern is repeated for five years, sampling year-round

Deleted:

Deleted: In order t

Deleted: Furthermore, in order t

Deleted: 2

Formatted: Font: 12 pt, Not Italic

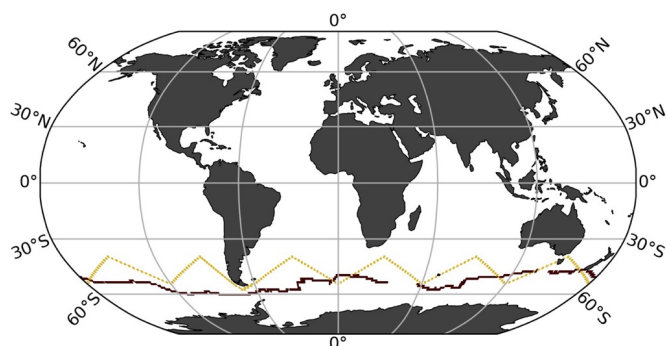
Formatted: Font: 12 pt

Formatted: Font: 12 pt, Not Italic

Deleted: Due to limited solar radiation that powers the Sairdron USVs, we let the sampling occur at a maximum latitude of 55° S.

Deleted: Considering the

314 ('Z_x10_5Y_YR'; 11,400 observations) and during southern hemisphere winter months only
 315 ('Z_x10_5Y_W'; 3,800 observations).



316
 317 **Figure 2:** SAILDRONE Uncrewed Surface Vehicle (USV) tracks representing the first circumnavigation around
 318 Antarctica from 2019 in maroon ('one-latitude' track; Sutton et al., 2021) and an alternative virtual route with
 319 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 YR	Z_x4 10Y W	Z_x10 5Y YR	Z_x10 5Y W
Saildrone track	NA	One-lat	One-lat	One-lat	One-lat	One-lat	One-lat	Zigzag	Zigzag	Zigzag	Zigzag	Zigzag
Years of sampling	NA	5	10	10	10	10	5	10	10	5	5	5
Duration of sampling	NA	Jan-Aug	Jan-Aug	Year-round	Jan-Aug	SO winter	SO winter	Year-round	SO winter	Year-round	Year-round	SO winter
Additional observations	NA	2,075	4,150	4,150	44,250	25,395	5,022	7,600	2,500	11,400		3,800
Global coverage increase (%)	NA	0.01	0.02	0.02	0.3	0.1	0.03	0.04	0.01	0.07		0.02
Mean bias (µatm)												
<i>Testbed period (1982-2016)</i>												
Globally	0.63	0.59	0.59	0.52	0.53	0.39	0.57	0.51	0.51	0.45	0.44	0.44
NORTH (35°N-90°N)	0.11	0.24	0.20	0.25	0.20	0.17	0.16	0.16	0.16	0.12	0.20	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	0.13	0.18	0.18
SOUTH (90°S-35°S)	1.4	1.3	1.2	1.1	1.1	0.80	1.2	1.1	1.1	1.0	0.87	0.87
SO winter months (JJA)	1.3	1.2	1.2	1.1	1.1	0.90	1.2	0.93	1.0	0.94	0.95	0.95
SO summer months (DJF)	0.070	0.11	0.15	0.10	0.15	0.019	0.11	0.25	0.073	0.16	0.066	0.066
<i>2006-2012-2016</i>												
Globally	0.51*	0.27	0.34	0.28	0.19	0.03	0.21	0.23	0.24	0.17	0.07	0.07
SOUTH (90°S-35°S)	1.6*	0.93	1.1	1.0	0.72	0.37	0.73	0.89	0.92	0.67	0.55	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	1.2	2.0	2.0
Mean RMSE (µatm)												
<i>Testbed period (1982-2016)</i>												
Globally	11.8	11.7	11.8	11.7	11.7	11.6	11.7	11.5	11.6	11.5	11.6	11.6
NORTH (35°N-90°N)	13.0	13.0	13.0	13.0	13.0	13.0	13.1	13.0	13.0	13.0	13.0	13.0
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	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	10.6	11.0	11.0
<i>2006-2012-2016</i>												
Globally	11.6*	11.6	11.4	11.3	11.3	11.2	11.6	11.0	11.2	11.1	11.4	11.4
SOUTH (90°S-35°S)	11.4*	11.1	11.0	10.7	10.6	10.4	10.9	10.0	10.6	9.7	10.6	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	9.6	10.3	10.3

320
 321 **Table 1.** Overview of the different sampling experiments tested in this study, and mean bias and RMSE (in µatm) for
 322 various time periods, latitude bands for all runs. **Bold values represent the best score for each category.** 'One-lat' =
 323 'one-latitude' track; incorporates the SAILDRONE USV route from Sutton et al. (2021). 'Zigzag' = potential meridional
 324 sampling. 'Additional observations' = number of 1°x1° monthly SAILDRONE USV observations in addition to SOCAT.
 325 J-A= January-August. YR = year-round. W = southern hemisphere winter. x4, x5, x10 and x13 = four, five, ten and
 326 13 USVs. SO winter = Southern Ocean winter months, i.e., June, July, August and also including September. *Average
 327 value of the mean of 2006-2016 and 2012-2016. The global coverage increase was calculated based on the total
 328 number of available 1982-2016 monthly 1°x1° observations from SOCAT (262,204 observations) and the Large
 329 Ensemble Testbed (17,290,470 observations).

330
 331 **2.5 Air-sea CO₂ flux**

Run name	5Y J-A	10Y J-A	10Y YR	x1
Saildrone track	One-lat	One-lat	One-lat	
Years of sampling	5	10	10	
# of SAILDRONES	1	1	1	
Duration of sampling	Jan-Aug	Jan-Aug	Year-round	
Total observations	2,075	4,150	4,150	
Global coverage increase (%)	0.01	0.02	0.02	

- Deleted:
- Deleted: SAILDRONE USV sampling patterns
- Formatted: Font: 10 pt
- Deleted: using the XGBoost Machine Learning algorithm (Gloege et al., 2021; Bennington et al., 2022a) to estimate surface ocean pCO₂
- Deleted: The 'one-latitude' (
- Deleted: o
- Formatted: Font: 10 pt
- Deleted:) track
- Deleted: ,
- Deleted: while the
- Deleted: z
- Deleted: track represents
- Deleted: future
- Deleted: (see Fig. 2)
- Deleted: The total
- Deleted: number of USV
- Deleted: (in bold)
- Deleted: represent
- Formatted: Pattern: Clear
- Formatted: Pattern: Clear
- Formatted: Pattern: Clear
- Formatted: Pattern: Clear
- Formatted: Pattern: Clear
- Formatted: Pattern: Clear
- Deleted: Note that all runs also included SOCAT coverage.

351 To assess the global ocean carbon sink associated with our pCO₂ reconstructions, air-sea CO₂
352 exchange was calculated for 1985 onward. Here, we computed air-sea CO₂ fluxes using the bulk
353 formulation with python package Seaflux.1.3.1 (<https://github.com/lukegre/SeaFlux>; Gregor et al.
354 2021; Fay et al., 2021). We calculated global and Southern Ocean flux in the same manner for 1)
355 the testbed ‘model truth’, 2) the ‘SOCAT-baseline’ and 3) the 10 experimental USV runs.

Deleted:

356 The net sea–air CO₂ flux was estimated using:

$$357 \text{ Flux} = k_w \cdot \text{sol} \cdot (\text{pCO}_2^{\text{ocn}} - \text{pCO}_2^{\text{atm}}) \cdot (1 - \text{ice})$$

358 where ‘k_w’ is the gas transfer velocity, ‘sol’ is the solubility of CO₂ in seawater (in units of mol
359 m⁻³ μatm⁻¹), ‘pCO₂^{ocn}’ is the partial pressure of surface ocean carbon (in μatm), either from the
360 ‘model truth’ or from the reconstructions, and pCO₂^{atm} (in μatm) is the partial pressure of
361 atmospheric CO₂ in the marine boundary layer. For GFDL, we used direct model output of
362 pCO₂^{atm}, while for CESM and CanESM2, pCO₂^{atm} was calculated individually, as the product of
363 surface xCO₂ and sea level pressure (the contribution of water vapor pressure was corrected for in
364 CESM and GFDL). Finally, to account for the seasonal ice cover in high latitudes, the fluxes were
365 weighted by 1 minus the ice fraction (‘ice’), i.e., the open ocean fraction. Inputs to the calculation
366 include EN4.2.2 salinity (Good et al., 2013), SST and ice fraction from NOAA Optimum
367 Interpolation Sea Surface Temperature V2 (OISSTv2) (Reynolds et al., 2002), and surface winds
368 and associated wind scaling factor from the European Centre for Medium-Range Weather
369 Forecasts (ECMWF ERA5 sea level pressure (Hersbach et al., 2020). Results presented show the
370 global and Southern Ocean (< 35° S) fluxes in units of Pg C yr⁻¹.

Deleted: pCO₂^{atm} from CESM was corrected for

371 Note that, reconstructions of pCO₂ for the ‘SOCAT-baseline’ and the experimental USV
372 runs are limited in their spatial extent to the open ocean (see Sect. 2.1; excluding coastal areas, the
373 Arctic Ocean and marginal seas). The same mask was thus also applied when calculating the flux
374 of the ‘model truth’, prior to comparison with the reconstructions.

Deleted:

375 3. Results

376 3.1 Performance metrics for the ‘SOCAT-baseline’ reconstruction

Deleted:

377 The mean bias for the entire testbed period (i.e., 1982-2016) is 0.63 μatm globally (Fig. 3a) and
378 1.4 μatm for the Southern Ocean (< 35° S; Table 1). Bias is much closer to zero for the mid-

Deleted: S

384 latitudes (between 35° S and 35° N; 0.23 μatm) and northern latitudes ($> 35^\circ$ N; 0.11 μatm) (**Fig.**
385 **3a**). There is a significant difference in bias considering southern hemisphere winter months (June,
386 July, August) versus summer months (December, January, February), with a global mean bias (for
387 1982-2016) of 1.3 μatm compared to 0.07 μatm , respectively (**Table 1**), due to the sparseness of
388 SOCAT observations from the southern hemisphere during the harsh winter season (**Fig. S5a**).
389 The mean RMSE for the entire testbed period (i.e., 1982-2016) is 11.8 μatm globally (**Fig. 3b**) and
390 11.5 μatm for the Southern Ocean (**Table 1**). RMSE is highest in the Eastern Tropical and
391 Southeastern Pacific Ocean and in the Southern Ocean, where the algorithm generally
392 overestimates pCO_2 (i.e., positive bias; **Fig. 3a**), with some exceptions in the Atlantic section. This
393 is consistent with the areas significantly undersampled by SOCAT (**Fig. S5b**). Except for these
394 areas, RMSE and bias is generally low (close to zero) in the open ocean, but show higher values
395 along coastlines (**Fig. 3b**).

Deleted: S

Deleted: 3

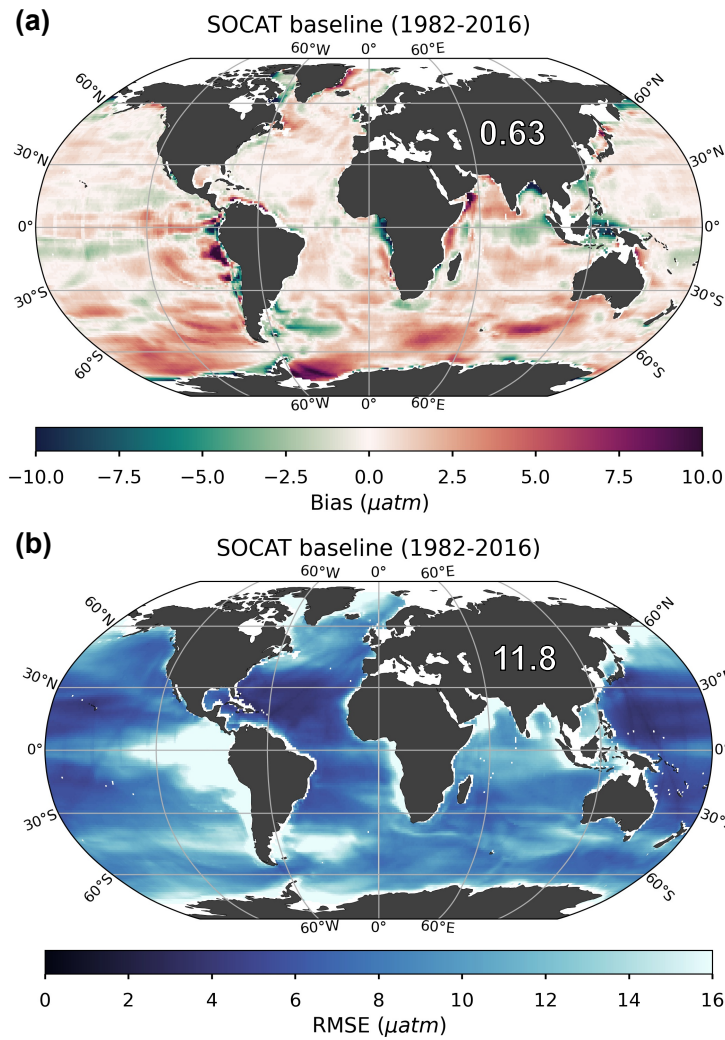
Deleted: 7

Deleted: 9

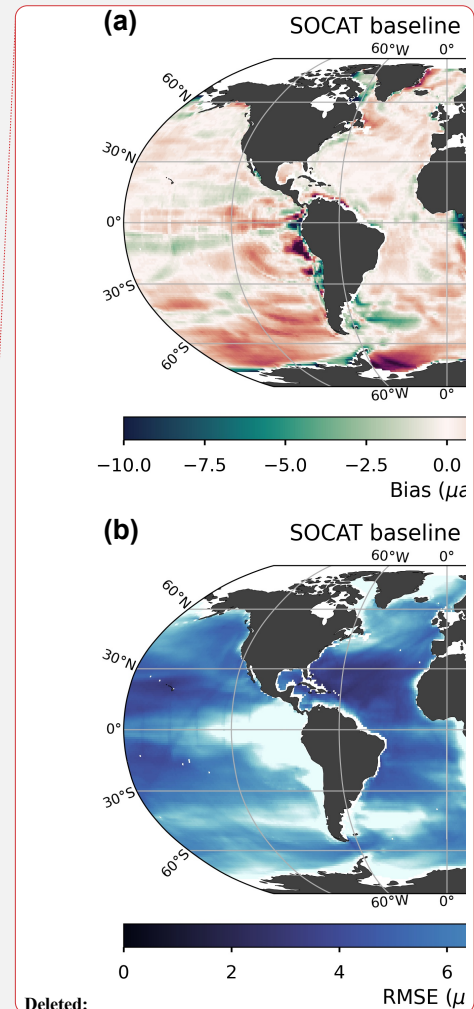
Deleted: 8

Deleted: S

Deleted: 3



403
 404 **Figure 3: Bias (a) and root-mean-squared error (RMSE) (b) for the ‘SOCAT-baseline’ (i.e., no USV) over the period**
 405 **of 1982 through 2016. The global mean bias and RMSE is 0.63 μatm and 11.8 μatm , respectively. Note that only the**
 406 **open ocean was considered in the reconstruction, so several areas were masked out prior to algorithm processing, such**
 407 **as the Arctic Ocean, coastal areas and marginal seas (no data; white areas in figures).**



Deleted:

Deleted: B

Deleted: when comparing the baseline machine learning reconstruction with the testbed ‘model truth’, averaged over the 75 ensemble members for the period of 1982 through 2016. The testbed was sampled based on SOCAT observations only (i.e., no USV).

Deleted: 7

Deleted: 9

Deleted: Red and green areas in **a** indicate regions where the reconstruction is biased high (i.e., overestimates pCO_2) and low (i.e., underestimates pCO_2), respectively. Generally, RMSE is highest in the East and South Pacific Ocean and in the Southern Ocean, where the algorithm also generally overestimates pCO_2 (positive bias; **a**).

425 3.2 Reconstruction improvements with Saildrone USV additions

426 Our presentation of global maps is limited to runs ‘x5_5Y_W’ (5,022 **monthly 1°x1°** observations)
427 and ‘Z_x4_10Y_YR’ (7,600 **monthly 1°x1°** observations). These runs were selected as they
428 represent observational schemes that are realistic in the near-term future considering logistics and
429 cost level, both non-meridional and meridional sampling, and different approaches to observing
430 duration and seasonal coverage. For the remaining runs, equivalent maps can be found in the
431 **Supplement**.

432 3.2.1 Bias

433 All Saildrone USV runs show a reduction in bias compared to the global mean 1982-2016
434 ‘SOCAT-baseline’ (Figs. 4a, S6). The improvement in bias is mainly due to lower reconstructed
435 pCO₂ values at southern latitudes, where the ‘SOCAT-baseline’ reconstruction generally
436 overestimates pCO₂ (Fig. 3a). The global mean bias for ‘zigzag’ run ‘Z_x4_10Y_YR’ is 0.51
437 μatm, a higher improvement (19 %) over the ‘SOCAT-baseline’ compared to the ‘one-latitude’
438 run ‘x5_5Y_W’ (11 % **mean** improvement; mean bias = 0.57 μatm;) (Fig. 4a; Table 1). Generally,
439 the ‘zigzag’ runs show higher improvements from the ‘SOCAT-baseline’ (19-31 % improvement;
440 **resulting** mean bias = 0.44-0.51 μatm) compared to the ‘one-latitude’ runs (7-19 % improvement;
441 **resulting** mean bias = 0.52-0.59 μatm) (Fig. S6; Table 1). However, the ‘one-latitude’-run
442 ‘x13_10Y_W’ that samples southern hemisphere winter months only, stands out with the lowest
443 global mean bias of 0.39 μatm, representing a 39 % **mean** improvement from the ‘SOCAT-
444 baseline’, **as well as reduced spread across the 75 ensemble members** (Table 1; Fig. S6; S8). This
445 run, however, has three or five times more observations (25,395) than ‘Z_x4_10Y_YR’ and
446 ‘x5_5Y_W’, respectively.

447 Compared to the entire testbed period, even larger improvements in global mean bias are
448 shown for the period of Saildrone USV additions (2006-2016 and 2012-2016; Figs. 4a vs. 4b,
449 Figs. S6 vs. S7). Compared to the ‘SOCAT-baseline’, run ‘x13_10Y_W’ results in a **mean** bias
450 improvement of 95 %, while the remaining ‘one-latitude’ runs and the ‘zigzag’ runs show **mean**
451 improvements up to 63 % and 85 %, respectively (Fig. S7).

452 Perhaps surprisingly, there is not a strong connection between the global or Southern Ocean
453 mean bias and the number of added USV observations (Fig. 5). The ‘one-latitude’ ‘high-sampling’

Deleted:

Deleted: 4

Deleted:

Deleted: S

Deleted:

Deleted: 4

Deleted: S

Deleted:

Deleted: S

Deleted: 4

Formatted: Font: Not Bold

Formatted: Highlight

Deleted: 4

Deleted: 5

Deleted:

Deleted: 5

468 run 'x13_10Y_J-A' (44,250 observations) show similar mean bias or is outperformed by all
469 'zigzag' runs as well as the 'one-latitude'-runs that restrict sampling to southern hemisphere winter
470 months (i.e., 'x5_5Y_W' and 'x13_10Y_W').

471 Considering the change in bias from year-to-year, the 'SOCAT-baseline' shows positive
472 bias at all latitudes in the beginning of the testbed period, before improvement occurs around 1990
473 (Fig. 6a). This is consistent with increasing SOCAT sampling with time for the period considered
474 here (i.e., up to 2016; Fig. S5c). As SOCAT observations are biased towards the northern
475 hemisphere (Fig. S5a, b), bias in the Southern Ocean (< 35° S) increases significantly starting in
476 the 2000s and remains high until the end of the testbed period (Fig. 6a). By adding USV sampling,
477 bias in the Southern Ocean improves over the 'SOCAT-baseline' around year 2000 (Fig. 6b-d;
478 Fig. S9), up to 6-12 years before to the introduction of additional samples in either 2006 or 2012.
479 This improvement is shown for the majority of the 75 ensemble members (Fig. S10). Run
480 'Z_x10_5Y_W', which has the lowest mean bias out of the 'zigzag' runs (Fig. 5), shows
481 improvement even further back in time, until the beginning of the testbed period (Fig. S9). While
482 the annual mean bias of the 'zigzag' runs varies rather consistently, there is a larger spread across
483 the 'one-latitude' runs (Fig. 6d).

Deleted:

Deleted: year

Deleted: time

Deleted: 3

Deleted: 3

Deleted:

Deleted: 6

Deleted: prior

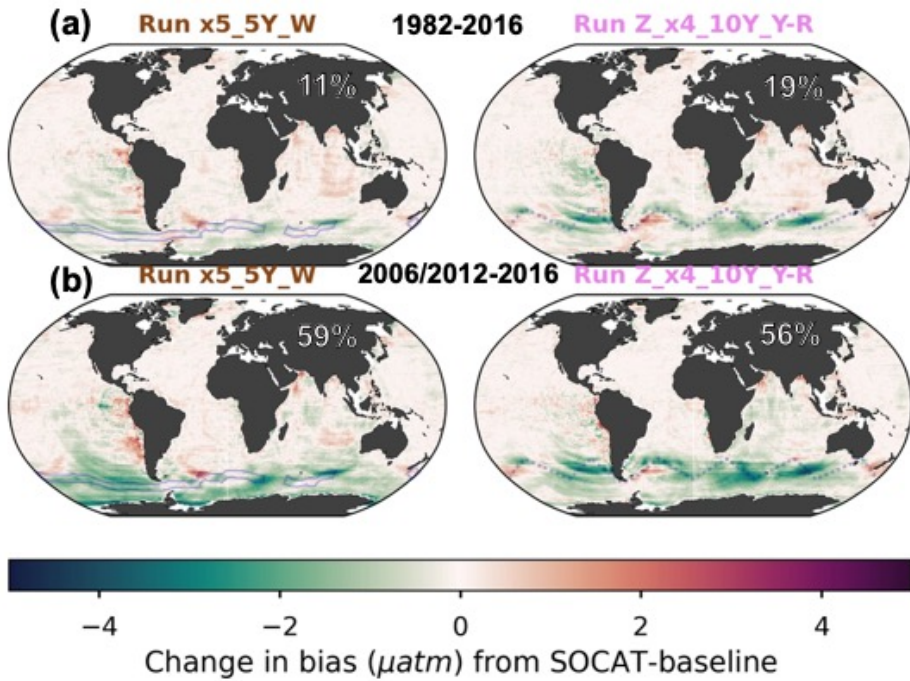
Formatted: Font: Bold

Deleted: 6

Deleted: vary

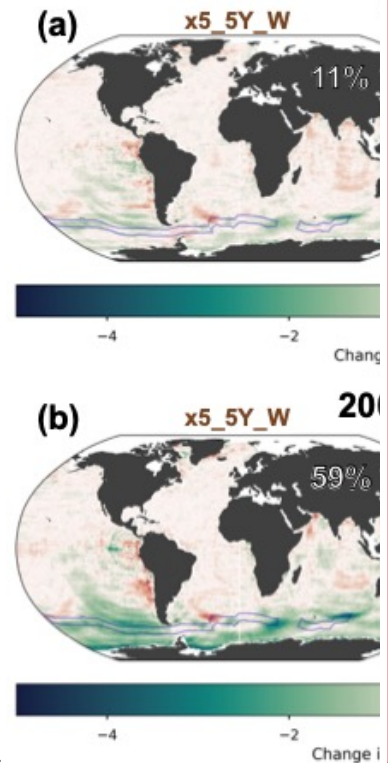
Deleted: in a similar manner

Deleted: between



496
497
498
499
500

Figure 4: Change in bias when comparing run ‘x5 5Y W’ and ‘Z x4 10Y YR’ to the ‘SOCAT baseline’ reconstruction, averaged over the duration of the testbed period (a; 1982-2016) and the period of USV additions (b; 2006-2012 or 2012-2016). The percent global improvement in absolute bias is shown on each panel.

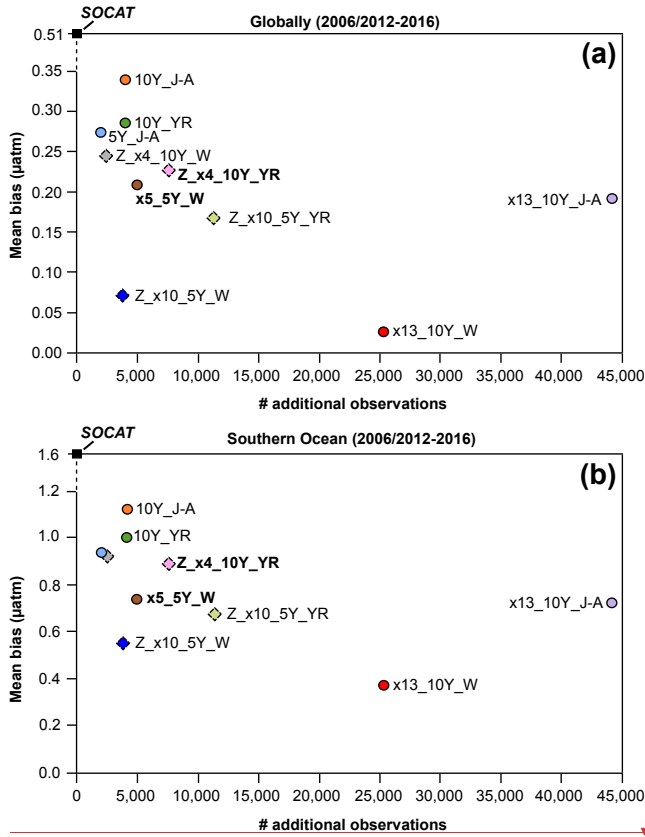


Deleted:

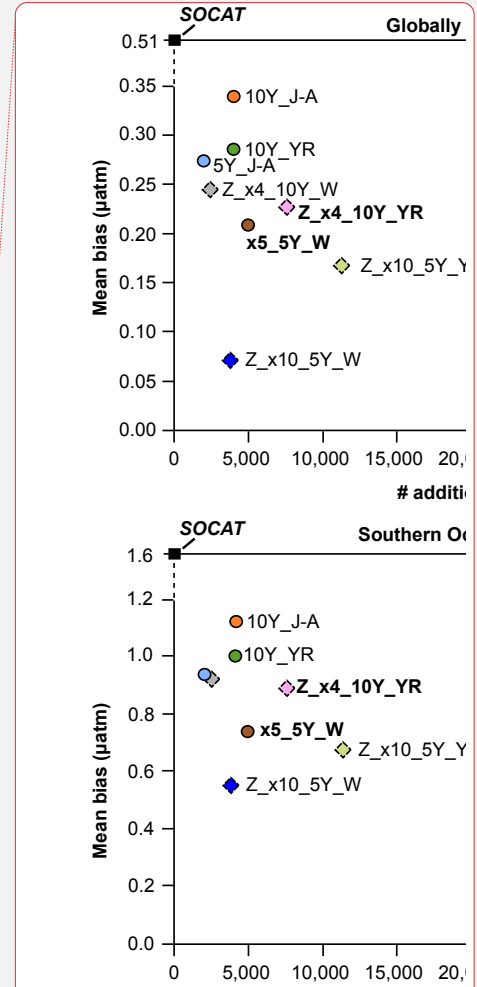
Deleted:

Deleted: Negative change in bias is found across the southern latitudes, indicating an improvement compared to the SOCAT baseline that overestimates pCO₂ (Figure 3a).

Deleted: Note that improvement is greater in the period of Saildrone USV additions compared to the entire testbed period. ...

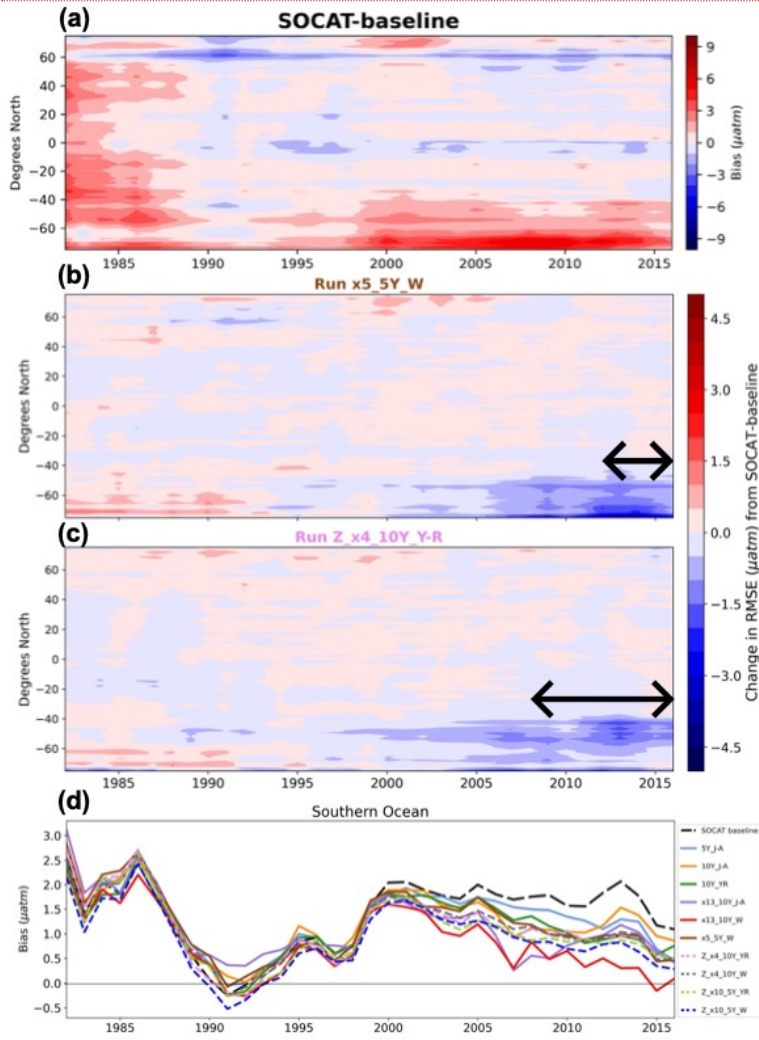


509
 510 **Figure 5:** Mean bias globally (a) and for the Southern Ocean (b) for the duration of Saildron USV sampling (2006-
 511 2016 or 2012-2016) for all runs presented in **Table 1**. Circles represent runs using the ‘one-latitude’ track, while
 512 diamonds represent ‘zigzag’ runs. Runs highlighted in bold correspond to the two selected runs mapped in **Figure 4**,
 513 **6**, **7** and **9**. Global (0.51 μatm) and Southern Ocean (1.6 μatm) bias values shown for the ‘SOCAT baseline’ (black
 514 squares) represent a mean of values for 2006-2016 (global = 0.52 μatm , S. Ocean = 1.63 μatm) and 2012-2016 (global
 515 = 0.51 μatm , S. Ocean = 1.56 μatm). ‘# additional observations’ = number of monthly $1^\circ \times 1^\circ$ USV observations in
 516 addition to SOCAT. Box plots illustrating the spread across the 75 ensemble members are shown in **Fig. S8**.



- Deleted: # additi
- Deleted: (Sutton et al., 2021)
- Deleted:
- Deleted: The SOCAT baseline run included 261,733 monthly $1^\circ \times 1^\circ$ observations.
- Formatted: Font: Bold

522

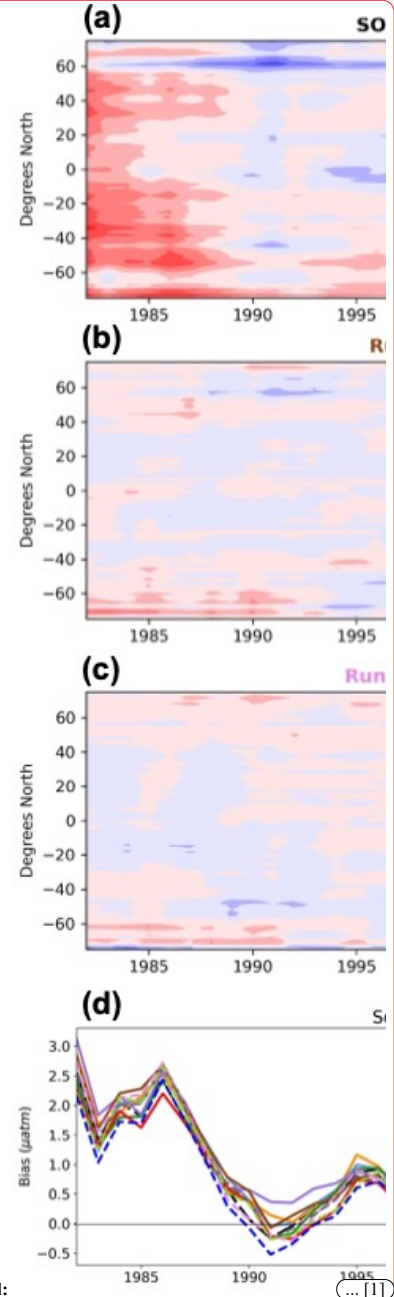


523

524 **Figure 6:** Zonal mean, annual mean Hovmöller of bias for the 'SOCAT-baseline' (a). Change in bias for run
 525 'x5_5Y_W' (b) and 'Z_x4_10Y_YR' (c) compared to the 'SOCAT-baseline' shown in (a). Improvement in bias in
 526 the Southern Ocean expands back in time well beyond the duration of USV additions for both runs (shown by arrows
 527 on each panel). Annual mean bias for the Southern Ocean (> 35° S) for all runs (d).

528

Deleted: Overall, there is not a strong correlation between bias and the number of observations, or duration of sampling.



Deleted:

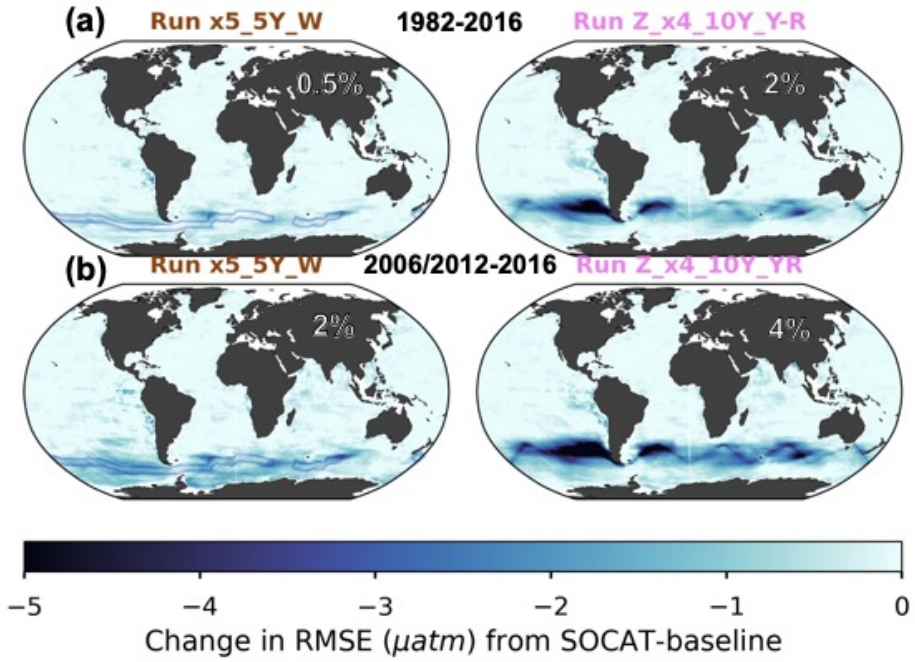
... 11

568 3.2.2 Root-mean squared error (RMSE)

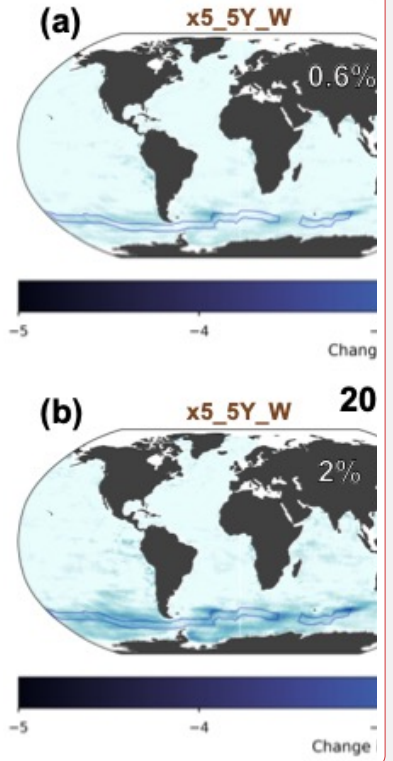
569 Similar to bias, improvements in RMSE are most significant during the period of USV additions
570 and within the Southern Ocean (Fig. 7a vs. 7b). For the duration of USV additions, the ‘one-
571 latitude’ runs show improvements in global mean RMSE of 1-3% (0.1-1% for 1982-2016), while
572 the ‘zigzag’ runs show higher improvements between 2-5% (1-3% for 1982-2016) (Figs. 7, S11,
573 S12). Mean RMSE is further reduced in the Southern Ocean by up to 16%, and during southern
574 hemisphere winter months (JJA) up to 21% (run ‘Z x10 5Y YR’; mean RMSE of 9.6 μatm ;
575 Table 1). There is minimal change in RMSE (or bias) during southern hemisphere summer months
576 (DJF; Fig. S13). The two ‘zigzag’ runs sampling year-round (‘Z x4 10Y YR’ and
577 ‘Z x10 5Y YR’) have the lowest RMSE values both globally and in the Southern Ocean (Fig. 8).
578 The spread across the 75 testbed members for each experiment is shown in Figure S14.

579 The ‘zigzag’ runs, as well as the ‘high-sampling’ ‘one-latitude’-runs (i.e., ‘x13_10Y_J-A’
580 and ‘x13_10Y_W’), show improvements compared to the ‘SOCAT-baseline’ from the initiation
581 of sampling (Figs. 9, S15, S16). The year-round ‘zigzag’ runs, however, show improvement in the
582 Southern Ocean from the beginning of the testbed period (Figs. 9c, d, S15). RMSE improvements
583 back in time are greater for all runs in the southern hemisphere winter months (Fig. S17).

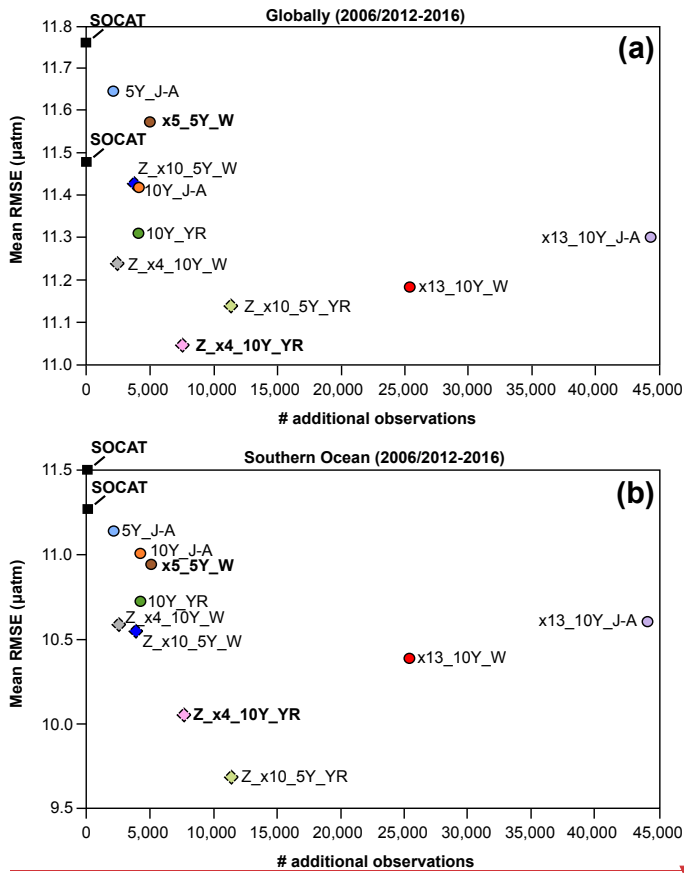
- Deleted: 4
- Deleted: 3
- Deleted: 2
- Deleted: 3
- Deleted: 8
- Deleted: 2
- Deleted: 7
- Deleted: 8
- Deleted: in the Southern Ocean by
- Deleted: 6
- Deleted: 6
- Deleted: 9
- Deleted: 5
- Deleted: 9
- Deleted: ’
- Formatted: Font: Bold
- Deleted:
- Deleted: 0
- Formatted: Font: Not Bold
- Deleted: 0
- Deleted: more significant
- Deleted: 1



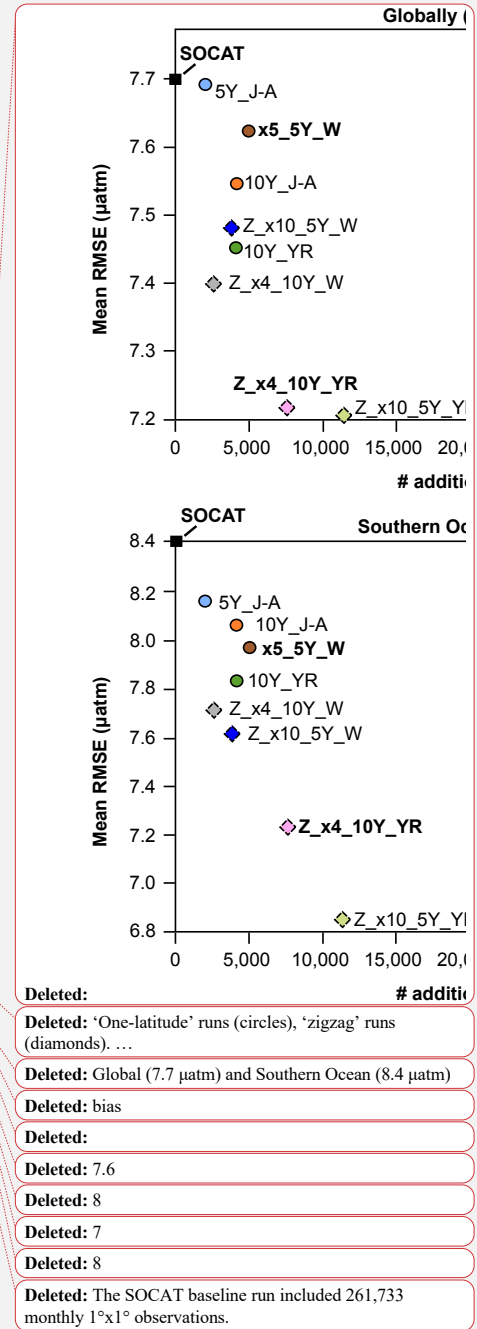
604
 605 **Figure 7:** Change in RMSE when comparing run 'x5 5Y W' and 'Z x4 10Y YR' to the 'SOCAT-baseline',
 606 averaged over the duration of the testbed period (a; 1982-2016) and the period of Saildrone USV additions (b; 2006-
 607 2012 or 2012-2016). The percent global improvement is shown on each panel.



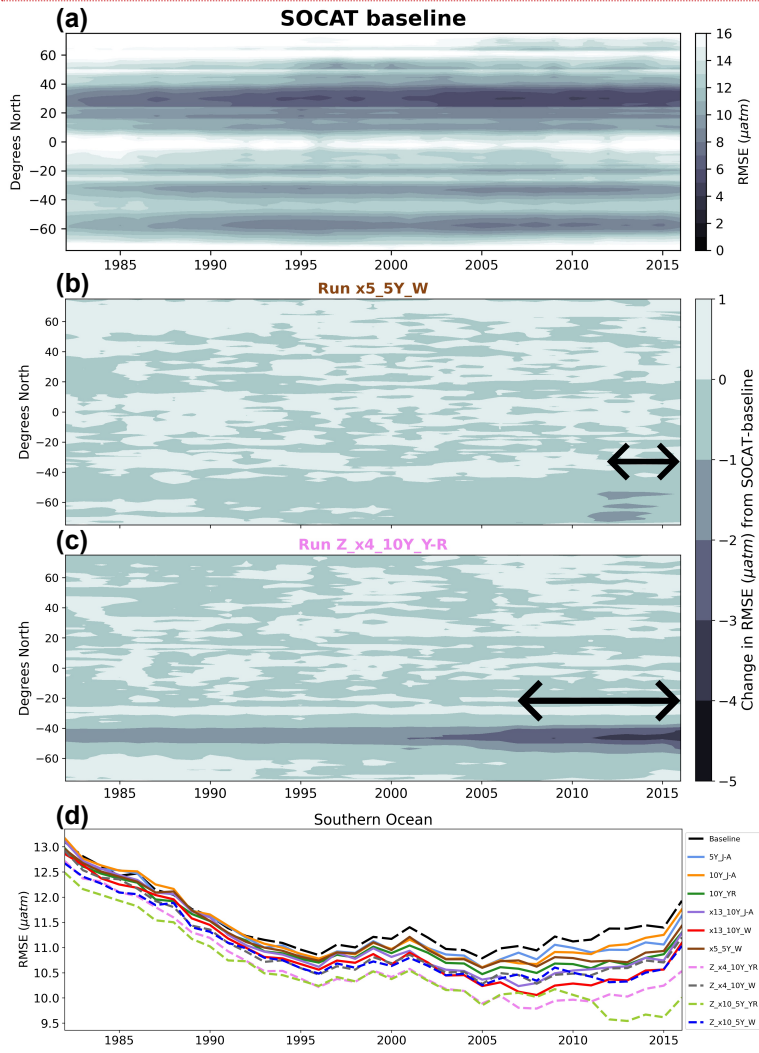
Deleted:
 Deleted:
 Deleted: reconstruction
 Deleted: Improvement in RMSE occurs mainly in southern latitudes (<35°S), where the baseline reconstruction shows high RMSEs (Fig. 3b).
 Deleted: Note the greater improvement for the period of USV additions compared to the entire testbed period.



616
 617 **Fig. 8:** Mean RMSE globally (a) and for the Southern Ocean (< 35° S; b) for the duration of Saildrone USV sampling
 618 (2006-2016 or 2012-2016) for all runs presented in Table 1. Circles represent runs using the 'one-latitude' track, while
 619 diamonds represent 'zigzag' runs. Runs highlighted in bold correspond to the two selected runs mapped in Figure 4,
 620 6, 7 and 9. RMSE values shown for the 'SOCAT_baseline' (black squares) represent a mean of values for 2006-2016
 621 (global = 11.5 μatm , S. Ocean = 11.3 μatm) and 2012-2016 (global = 11.8 μatm , S. Ocean = 11.5 μatm). '# additional
 622 observations' = number of monthly 1°x1° USV observations in addition to SOCAT. Box plots illustrating the spread
 623 across the 75 ensemble members are shown in Fig. S14.



636

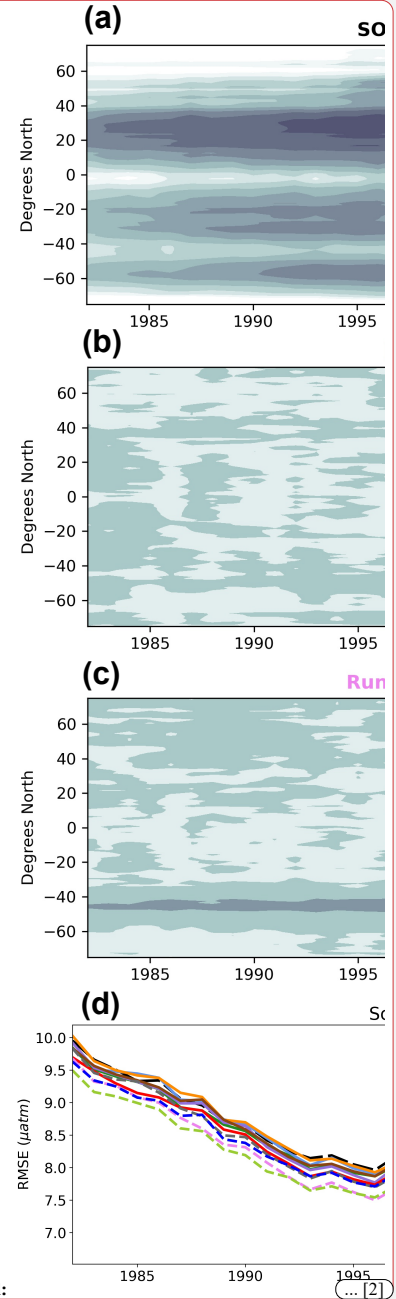


637

638 **Figure 9:** Zonal mean, annual mean Hovmöller of RMSE for the 'SOCAT-baseline' (a). Change in RMSE for run
 639 'x5_5Y_W' (b) and 'Z x4_10Y_YR' (c) compared to the 'SOCAT-baseline'. Run 'Z x4_10Y_YR' shows
 640 improvement in RMSE within the Southern Ocean, which expand well beyond the duration of Salldrone USV
 641 additions (shown by arrow on panel). Annual mean RMSE for the Southern Ocean (> 35° S) for all runs (d).

642

Deleted: Overall, there is not a strong correlation between increasing number of observations or duration of sampling and decreasing RMSE.



Deleted:

681 3.3 Impact on the air-sea CO₂ flux with Saildrone USV additions

682 Air-sea flux was calculated in the same manner for both the ML reconstructions and the ‘model
683 truth’, which allows for the isolation of the impact of different sampling strategies, as mediated by
684 the pCO₂ reconstruction, on fluxes (see Sect. 2.5). These flux estimates are made to inform
685 understanding of the errors that may exist in CO₂ flux estimates derived from pCO₂
686 reconstructions, and how new sampling could address these errors. Flux estimates represent the
687 average of the 75 members of the LET in each case, and are not estimates of real-world fluxes.

688 Compared to the ‘model truth’, the ‘SOCAT-baseline’ reconstruction underestimates the
689 global and Southern Ocean sink by 0.11-0.13 Pg C yr⁻¹ over 1982-2016 (Fig. 10; Table S1).
690 Regardless of sampling pattern, adding Saildrone USV observations increases both the global and
691 Southern Ocean mean sink compared to the ‘SOCAT-baseline’ (Figs. 10, S18). The ‘one-latitude’
692 runs show an increase of 0.01-0.03 Pg C yr⁻¹ (2-6 % strengthening) of the Southern Ocean sink
693 (1982-2016), while the ‘zigzag’ runs lead to an even stronger sink by 0.04-0.06 Pg C yr⁻¹ (7-11 %
694 strengthening) (Table S2). When averaging over the years of Saildrone USV sampling addition
695 (i.e., 2006-2012 and 2012-2016), the Southern Ocean sink increases up to 0.09 Pg C yr⁻¹ (14 %
696 strengthening) for the ‘one-latitude’ runs and up to 0.1 Pg C yr⁻¹ (15 % strengthening) for the
697 ‘zigzag’ runs (Table S2). These same features are found for the global ocean (Fig. S18; Table
698 S2).

699 All of the ‘zigzag’ runs quite closely match both the global and Southern Ocean ‘model
700 truth’ air-sea CO₂ flux for the duration of sample additions (Figs. 10, S18). Except for the first
701 couple of years of sample addition for the ‘high-sampling’-run ‘x13_10Y_J-A’, none of the ‘one-
702 latitude’ runs can match the ‘model truth’ air-sea CO₂ flux, instead they all underestimate the flux
703 (Figs. 10, S18). The ‘zigzag’ runs have impact on the air-sea flux from an earlier date, starting to
704 pull the results away from the ‘SOCAT-baseline’ and toward the ‘model truth’ already in the late-
705 1990s, while the ‘one-latitude’ runs do the same about a decade later (Figs. 10, S18).

Deleted: direct comparison of the differences in

Formatted: Subscript

Deleted: These fluxes

Deleted:

Deleted: 2

Deleted:

Deleted: 2

Deleted: 3

Deleted: 3

Deleted: 2

Deleted: 3

Deleted: 2

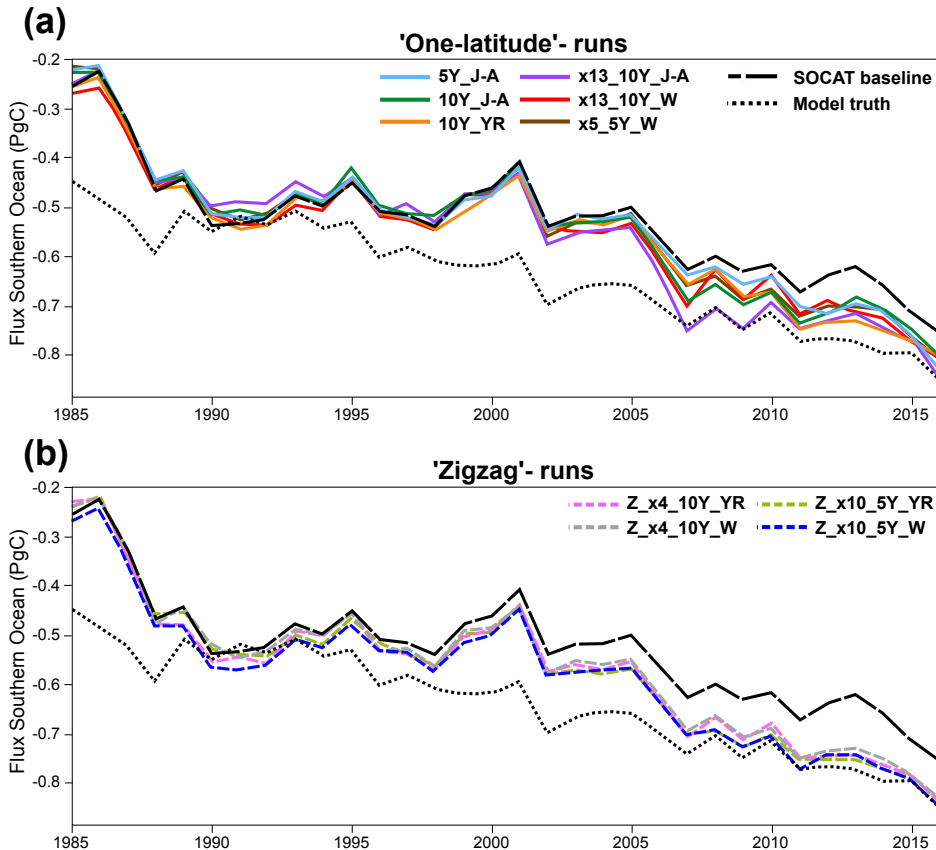
Deleted: are able to

Deleted: as

Deleted: 2

Deleted:

Deleted: 2



722
723 **Figure 10:** Southern Ocean (< 35° S) annually averaged air-sea CO₂ flux for the 'SOCAT baseline' (black dashed
724 line), 'model truth' (black dotted line) 'one-latitude' runs (a; solid lines) and 'zigzag' runs (b; dashed lines).

725
726
727 **4. Discussion**

728 We have tested the pCO₂-Residual reconstruction method with the Large Ensemble Testbed (LET)
729 to estimate its fidelity and understand how new samples could increase skill. We find that,
730 regardless of the chosen Sailability USV sampling pattern, the reduction in **mean bias** and **mean**
731 **RMSE** compared to the 'SOCAT baseline' is most prominent within the Southern Ocean (< 35°
732 S) during the period of which Sailability USV observations were added (Figs. 4, 6, 7, 9). However,
733 it is important to mention that additional Southern Ocean sampling also improves pCO₂

Deleted:

Deleted: , averaged over the 75 ensemble members

Deleted: Compared to the SOCAT baseline, regardless of sampling pattern, the Sailability USV additions lead to an increased ocean sink. The 'zigzag' runs generate a stronger sink compared to the 'one-latitude' runs, and closely match the 'model truth' for the duration of sample additions.

Deleted: both

Deleted:

743 reconstructions globally (Figs. 5a, 8a). Based on our experiments, a combination of factors
 744 improve global and Southern Ocean pCO₂ reconstructions, including the type of sampling pattern
 745 and seasonality of sampling, and to some extent, the number of additional observations.
 746 Importantly, increasing the number of observations or duration of sampling (5 vs. 10 years) is not
 747 the sole determining factor for improving the reconstructions (Figs. 5, 8). This is best demonstrated
 748 by the ‘high-sampling’-run ‘x13_10Y_J-A’ (44,250 observations), which does not provide
 749 significantly better reconstructions, or is even outperformed, by runs with 2-18 times fewer
 750 observations. The runs that produce lower mean RMSE do include data throughout southern
 751 hemisphere winter (Figs. 8, 9d). Run ‘x13_10Y_J-A’ does not include more than a few
 752 observations in the month of August, as it follows the temporal pattern of the real-world ‘one-
 753 latitude’ Saildrone USV expedition (Fig. S2; Sutton et al., 2021). The ‘one-latitude’ runs ‘10Y_J-
 754 A’ and ‘10Y_YR’ are directly comparable in terms of sample duration, spatial extent and number
 755 of observations (Table 1), but the latter, which covers all months, always shows lower mean
 756 RMSE and bias (Figs. 5, 6d, 8, 9d). These examples attest to the importance of addressing the
 757 issue of significant undersampling in the Southern Ocean during the winter season (Figs. S5a, b).

758 Another important comparison is the ‘one-latitude’-run ‘x5_5Y_W’ (5,022 observations)
 759 and ‘zigzag’-run ‘Z_x10_5Y_W’ (3,800 observations) that both sample during southern
 760 hemisphere winter months over a five-year period (Table 1), where the ‘zigzag’-run consistently
 761 performs better even though it includes fewer observations (Figs. 5, 8). Most of the runs that
 762 perform similar to, or outperform, the above-mentioned ‘high-sampling’-run ‘x13_10Y_J-A’
 763 (44,250 observations), sample in a ‘zigzag’ pattern. Out of all 10 runs, the ‘year-round’ ‘zigzag’
 764 runs (‘Z_x4_10Y_YR’ and ‘Z_x10_5Y_YR’) are most able to reduce the mean error as shown by
 765 the lowest RMSE values (Figs. 8, 9d). A recent study performed similar sampling experiments as
 766 shown here, by comparing sampling from different types of autonomous platforms to a ‘SOCAT-
 767 baseline’ (Djeutchouang et al., 2022). They emphasized the importance of capturing the significant
 768 differences in pCO₂ that exist across meridional gradients during summer and winter months (up
 769 to 15 µatm; Djeutchouang et al., 2022). The meridional coverage provided by the ‘zigzag’ runs
 770 could explain why these runs generally outperform the ‘one-latitude’ runs in our study, and show
 771 significant reduction in both RMSE and bias, even though the global pCO₂ data density is raised
 772 by as little as 0.01-0.07%.

Deleted: seems to be important in order to

Deleted: both the

Deleted: and

Deleted: e

Deleted: mainly

Deleted: but also

Deleted: 2-18 times less

Deleted: , but that cover the full

Deleted: 5, 6d,

Deleted: 1

Deleted: 3

Deleted: magnitude of

Deleted:

Deleted: 4

787 The greatest reduction In mean bias out of all runs is shown by run 'x13_10Y_W' (Figs.
788 5, 6d), which represents 'one-latitude' 'high-sampling' (i.e., 25,395 observations) during southern
789 hemisphere winter months only. This sampling strategy seems thus to have a higher ability to
790 reduce the ML model's tendency to overestimate pCO₂ in the Southern Ocean compared to any of
791 the meridional ('zigzag') runs. However, it should be noted that run 'x13_10Y_W' covers areas
792 south of 55° S (Fig. S4), and its improvement in mean bias (and mean RMSE) is particularly
793 prevalent at these high latitudes (e.g., Figs. S7, S9, S12, S17). Whether or not this run is, in fact,
794 feasible with current or future technology is uncertain as parts of the southernmost tracks
795 potentially cover the Southern Ocean ice zone (Fig. S19), and solar radiation for solar-powered
796 platforms and sensors becomes very limited during winter south of 55° S. Furthermore, this
797 particular sampling strategy requires 13 USVs, and so would be the most costly of the observing
798 scenarios. Although run 'x13_10Y_W' demonstrates the highest reduction in mean bias out of all
799 runs, the 'zigzag' runs still reduce mean bias in the Southern Ocean by 44-65 % (vs. 77 % for run
800 'x13_10Y_W').

Deleted: i

Deleted: however

Deleted: 2

Deleted: such

Deleted: 5

Deleted: 6

Deleted: 8

Deleted: 0

Deleted: 3

Deleted: thus

801 Overall, the 'zigzag' runs include significantly fewer observations, require fewer USVs,
802 collect samples over the same duration, or even half the time as run 'x13_10Y_W', cover areas
803 north of 55°S and within the ice-free zone, and show major improvement in the reconstruction of
804 pCO₂, attested to by reductions in both bias and RMSE. The 'zigzag' runs also closely match both
805 the global and Southern Ocean 'model truth' air-sea CO₂ flux for the duration of sample additions
806 (Figs. 10, S18). It also appears that the 'zigzag' runs generally have a greater impact on both the
807 pCO₂ reconstruction and the air-sea flux further back in time, starting to deviate from the 'SOCAT-
808 baseline' earlier compared to the 'one-latitude' runs (Figs. 6, 9, 10, S9, S15, S17, S18). Even the
809 'zigzag' scenarios with the least number of USVs (e.g., 'Z_x4_10Y_YR') reduces Southern Ocean
810 reconstruction bias and RMSE by up to 46 % and 11%, respectively, and could provide a basis
811 for realistic future Southern Ocean pCO₂ sampling campaigns.

Deleted: less

Deleted: 2

Deleted:

Deleted: 6

Deleted: 0

Deleted: 1

Deleted: 2

Deleted: 3

812 The main motivation for improving surface ocean pCO₂ reconstructions is so that we can
813 more accurately estimate the current and future oceanic uptake of anthropogenic carbon. The
814 Southern Ocean is a significant carbon sink, but estimates of the air-sea CO₂ flux diverge
815 substantially in this region (Takahashi et al., 2009; Landschützer et al., 2014, 2015; Rödenbeck et
816 al., 2015; Williams et al., 2017; Gray et al., 2018; Gruber et al., 2019; Bushinsky et al., 2019; Long

835 et al., 2021; Fay and McKinley, 2021; Wu et al., 2022). Southern Ocean estimates incorporating
836 observations from biogeochemical floats have shown a significantly weaker sink compared to
837 those based only on observations from ships (Williams et al., 2017; Gray et al., 2018; Bushinsky
838 et al., 2019). Bushinsky et al. (2019) and Hauck et al. (2023) performed similar sampling
839 experiments as presented here, by comparing ML surface ocean pCO₂ reconstructions based on
840 SOCAT vs. additional SOCCOM or ideal virtual floats. These studies showed that SOCAT
841 sampling alone overestimates the CO₂ uptake in the Southern Ocean, and that additional floats
842 reduce this overestimation, leading to a decreased (weakened) ocean carbon sink. In contrast, we
843 find that the pCO₂-Residual method underestimates the CO₂ uptake with only SOCAT sampling,
844 and that adding USVs increased (strengthened) the Southern Ocean and global ocean sink by up
845 to 0.1 Pg C yr⁻¹ (Figs. 10, S18; Table S2).

846 Going forward, additional studies are needed to better understand why these results suggest
847 a different direction of the sink change with additional sampling. These differences could stem
848 from the use of different reconstruction methods assessed. Hauck et al. (2023) used the MPI-SOM-
849 FFN and CarboScope/Jena-MLS reconstruction methods, while we use the pCO₂-Residual
850 method. Another substantial difference between the studies is the models and numbers of ensemble
851 members used as the testbed. Hauck et al. (2023) use a single hindcast model, while we use 25
852 members each from three Earth System Models. We find substantial spread across these 75
853 members (Figs. S8, S10, S14, S16), indicating that model structure and internal variability
854 significantly impact results. Our study and Hauck et al. (2023) use different approaches for the
855 calculation of fluxes, which could also be a factor. Targeted, coordinated studies using multiple
856 reconstruction approaches with consistent testbed structures and experimental approaches are
857 clearly needed (Rödenbeck et al., 2015). Despite this need for this additional work, studies do
858 agree that additional Southern Ocean observations could significantly improve reconstructions of
859 air-sea CO₂ fluxes.

860 What else can we learn using the model testbed? The 'SOCAT-baseline' demonstrates a
861 weakening of the global and Southern Ocean carbon sink starting in the 1990s with a peak around
862 year 2000 (Figs. 10, S18), which is in broad agreement with various data products using real-world
863 SOCAT data (e.g., Gruber et al., 2019; Landschützer et al., 2015; Bushinsky et al., 2019;
864 Bennington et al., 2022; Gloege et al., 2022). Peaks in bias and RMSE coincide in time with the

Deleted: alone

Deleted: Southern Ocean

Deleted: They

Deleted: by

Deleted: ng

Deleted: the

Formatted: Subscript

Deleted: , the Southern Ocean carbon sink (mean of the period of float additions; 2015-2017) decreased (weakened) by 0.4 Pg C yr⁻¹. In contrast,

Deleted: by using a model testbed, we show that

Formatted: Subscript

Formatted: Subscript

Deleted: 2

Deleted: 3), which is a significant fraction of the uncertainty in the global ocean carbon sink (0.4 Pg C yr⁻¹; Friedlingstein et al., 2022

Deleted: Fed with real-world SOCAT data, the global mean air-sea flux estimate from the pCO₂-Residual method is similar to other available products (Bennington et al., 2022a), suggesting that other products may also underestimate the Southern Ocean carbon sink due to the spatio-temporal distribution of SOCAT data. Our experiments suggest that targeted USV observations could reduce this underestimation of the ocean carbon sink.

Formatted: Subscript

Formatted: Font: Bold, Not Highlight

Formatted: Not Highlight

Formatted: Font: Bold, Not Highlight

Formatted: Font: Bold

Formatted: Font: Bold

Formatted: Font: Bold

Formatted: Subscript

Deleted: ¶

Deleted:

Deleted: 2

890 weakening sink (Figs. 6d, 9d). As shown by Figure 10, this ‘low sink’ is significantly exaggerated
891 compared to the ‘model truth’. To better understand this discrepancy, we performed an additional
892 experiment based on run ‘Z_x10_5Y_YR’, but assumed sampling every year for the entire testbed
893 period (i.e., 1982-2016). There is now a significant reduction in the temporal variability of
894 reconstruction bias; with the additional 35-year USV sampling, the reconstructed Southern Ocean
895 air-sea CO₂ flux closely matches the ‘model truth’ for the entire testbed duration (Fig. S20). This
896 suggests that the large decadal variability of air-sea CO₂ fluxes since the 1980s, and the weak
897 anomaly in the Southern Ocean carbon sink in the early 2000s (Le Quéré et al., 2007; Landschützer
898 et al., 2015; Gruber et al., 2019; Bennington et al., 2022a,b; Friedlingstein et al., 2023), may be at
899 least partially attributable to undersampling of the Southern Ocean. This is in agreement with the
900 float sampling experiments performed by Hauck et al. (2023), attributing the strong decadal
901 variability to sparse and skewed SOCAT data distributions. We will further explore this issue in
902 future work. Still, this preliminary experiment suggests that interpretations of trends and variability
903 of the global and Southern Ocean carbon sink should be considered with caution.

904 5. Conclusions

905 By using the Large Ensemble Testbed (LET), we show that targeted meridional and winter
906 sampling in the Southern Ocean can improve global and Southern Ocean ML surface ocean pCO₂
907 reconstructions. Significant improvements are possible by raising the global pCO₂ data density by
908 as little as 0.01-0.07%. Further, we find that this modest amount of additional Sairdrone USV
909 sampling increases the global and Southern Ocean air-sea CO₂ flux by up to 0.1 Pg C yr⁻¹, a
910 quantity equivalent to 25 % of the uncertainty in the ocean carbon sink (0.4 Pg C yr⁻¹;
911 Friedlingstein et al., 2023). Our findings are consistent with previous studies suggesting that
912 additional observations during southern hemisphere winter months and covering meridional
913 gradients can reduce uncertainties and biases in the reconstructions (Lenton et al., 2006; Monteiro
914 et al., 2010; Djeutchouang et al., 2022; Mackay et al., 2022). As opposed to other autonomous
915 platform approaches, Sairdrone USVs obtain in situ pCO₂ observations with uncertainties
916 equivalent to the highest-quality observations collected by research ships ($\pm 2 \mu\text{atm}$; Sabine et al.,
917 2020; Sutton et al., 2021), and can operate at a high speed so that the spatial extent and seasonal
918 cycle of meridional gradients can be covered. The approach of combining high-accuracy Sairdrone
919 USV and SOCAT observations represents thus a promising solution to improve future surface

Deleted: The results from this experiment show

Deleted: 14

Deleted: 2

Deleted: 4

924 ocean pCO₂ reconstructions and the accuracy of the ocean carbon sink. Lastly, we show that the
925 large variability in bias, and the weakening of the global and Southern Ocean carbon sink in the
926 2000s, may be partially an artefact of Southern Ocean undersampling.

927 **Code availability**

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

929 **Data availability**

930 The Large Ensemble Testbed is publicly available at
931 https://figshare.com/collections/Large_ensemble_pCO2_testbed/4568555.

932

933 **Author contribution**

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

937 **Competing interests**

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

939 **Acknowledgements**

940 We acknowledge funding from NOAA through the Climate Observations and Monitoring Program
941 (Award #NA20OAR4310340) and from NSF through the LEAP STC (Award #2019625). This is
942 PMEL contribution 5549. We would also like to acknowledge and thank Val Bennington, Julius
943 Busecke, ~~Devan Samant~~ and ~~Abby Shaum~~ for providing technical support.

944

945 **References**

946

947 Bakker, D. C. E., Pfeil, B., Landa, C. S., Metzl, N., O'Brien, K. M., Olsen, A., Smith, K., Cosca,
948 C., Harasawa, S., Jones, S. D., Nakaoka, S., Nojiri, Y., Schuster, U., Steinhoff, T., Sweeney, C.,

Deleted: and

Deleted:

Formatted: Font: 12 pt

951 Takahashi, T., Tilbrook, B., Wada, C., Wanninkhof, R., Alin, S. R., Balestrini, C. F., Barbero, L.,
952 Bates, N. R., Bianchi, A. A., Bonou, F., Boutin, J., Bozec, Y., Burger, E. F., Cai, W.-J., Castle, R.
953 D., Chen, L., Chierici, M., Currie, K., Evans, W., Featherstone, C., Feely, R. A., Fransson, A.,
954 Goyet, C., Greenwood, N., Gregor, L., Hankin, S., Hardman-Mountford, N. J., Harlay, J., Hauck,
955 J., Hoppema, M., Humphreys, M. P., Hunt, C. W., Huss, B., Ibáñez, J. S. P., Johannessen, T.,
956 Keeling, R., Kitidis, V., Körtzinger, A., Kozyr, A., Krasakopoulou, E., Kuwata, A., Landschützer,
957 P., Lauvset, S. K., Lefèvre, N., Lo Monaco, C., Manke, A., Mathis, J. T., Merlivat, L., Millero, F.
958 J., Monteiro, P. M. S., Munro, D. R., Murata, A., Newberger, T., Omar, A. M., Ono, T., Paterson,
959 K., Pearce, D., Pierrot, D., Robbins, L. L., Saito, S., Salisbury, J., Schlitzer, R., Schneider, B.,
960 Schweitzer, R., Sieger, R., Skjelvan, I., Sullivan, K. F., Sutherland, S. C., Sutton, A. J., Tadokoro,
961 K., Telszewski, M., Tuma, M., van Heuven, S. M. A. C., Vandemark, D., Ward, B., Watson, A.
962 J., and Xu, S.: A multi-decade record of high-quality $f\text{CO}_2$ data in version 3 of the Surface Ocean
963 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)
964 [2016](https://doi.org/10.5194/essd-8-383-2016), 2016.

965 Bakker, D. C. E., Alin, S. R., Becker, M., Bittig, H. C., Castaño-Primo, R., Feely, R. A., Gkritzalis,
966 T., Kadono, K., Kozyr, A., Lauvset, S. K., Metzl, N., Munro, D. R., Nakaoka, S., Nojiri, Y., O'Brien,
967 K. M., Olsen, A., Pfeil, Benjamin, P., Denis, S., Tobias, S., Kevin F., Sutton, A. J., Sweeney, C.,
968 Tilbrook, B., Wada, C., Wanninkhof, R., Willstrand W. A., Akl, J., Apelthun, L. B., Bates, N.,
969 Beatty, C. M., Burger, E. F., Cai, W., Cosca, C. E., Corredor, J. E., Cronin, M., Cross, J. N., De
970 Carlo, E. H., DeGrandpre, M. D., Emerson, S. R., Enright, M. P., Enyo, K., Evans, W., Frangoulis,
971 C., Fransson, A., García-Ibáñez, M. I., Gehrung, M., Giannoudi, L., Glockzin, M., Hales, B.,
972 Howden, S. D., Hunt, C. W., Ibáñez, J. S. P., Jones, S. D., Kamb, L., Körtzinger, A., Landa, C.
973 S., Landschützer, P., Lefèvre, N., Lo Monaco, C., Macovei, V. A., Maenner J. S., Meinig, C.,
974 Millero, F. J., Monacci, N. M., Mordy, C., Morell, J. M., Murata, A., Musielewicz, S., Neill, .,
975 Newberger, T., Nomura, D., Ohman, M., Ono, T., Passmore, A., Petersen, W., Petihakis, G.,
976 Perivoliotis, L., Plueddemann, A. J., Rehder, G., Reynaud, T., Rodriguez, C., Ross, A. C.,
977 Rutgersson, A., Sabine, C. L., Salisbury, J. E., Schlitzer, R., Send, U., Skjelvan, I., Stamataki, N.,
978 Sutherland, S. C., Sweeney, C., Tadokoro, K., Tanhua, T., Telszewski, M., Trull, T., Vandemark,
979 D., van Ooijen, E., Voynova, Y. G., Wang, H., Weller, R. A., Whitehead, C., Wilson, D.: Surface
980 Ocean CO_2 Atlas Database Version 2022 (SOCATv2022) (NCEI Accession 0253659), NOAA

981 National Centers for Environmental Information [dataset], <https://doi.org/10.25921/1h9f-nb73>,
982 2022.

983 Bennington, V., Galjanic, T., and McKinley, G. A.: Explicit Physical Knowledge in Machine
984 Learning for Ocean Carbon Flux Reconstruction: The pCO₂-Residual Method, Journal of
985 Advances in Modeling Earth Systems, 14(10), <https://doi.org/10.1029/2021ms002960>, 2022a.

986 Bennington, V., Gloege, L., and McKinley, G. A.: Variability in the global ocean carbon sink from
987 1959 to 2020 by correcting models with observations, Geophysical Research Letters, 49(14),
988 <https://doi.org/10.1029/2022GL098632>, (2022b).

989 Bushinsky, S. M., Landschützer, P., Rödenbeck, C., Gray, A. R., Baker, D., Mazloff, M. R.,
990 Resplandy, L., Johnson, K. S., and Sarmiento, J. L.: Reassessing Southern Ocean air-sea CO₂ flux
991 estimates with the addition of biogeochemical float observations, Global Biogeochemical Cycles,
992 33(11), 1370-1388, <https://doi.org/10.1029/2019GB006176>, 2019.

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

996 Deser, C., Phillips, A., Bourdette, V., and Teng, H.: Uncertainty in climate change projections: the
997 role of internal variability, Climate Dynamics, 38, 527-546, <https://doi.org/10.1007/s00382-010-0977-x>, 2012

999 Djeutchouang, L. M., Chang, N., Gregor, L., Vichi, M., and Monteiro, P. M. S.: The sensitivity of
1000 pCO₂ reconstructions to sampling scales across a Southern Ocean sub-domain: a semi-idealized
1001 ocean sampling simulation approach, Biogeosciences, 19, 4171-4195, <https://doi.org/10.5194/bg-19-4171-2022>, 2022

1003 Fay, A. R., Lovenduski, N. S., McKinley, G. A., Munro, D. R., Sweeney, C., Gray, A. R.,
1004 Landschützer, P., Stephens, B. B., Takahashi, T., and Williams, N.: Utilizing the Drake Passage
1005 Time-series to understand variability and change in subpolar Southern Ocean pCO₂,
1006 Biogeosciences, 15(12), 3841-3855, <https://doi.org/10.5194/bg-15-3841-2018>, 2018.

1007 Fay, A. R., and McKinley, G. A.: Observed regional fluxes to constrain modeled estimates of the
1008 ocean carbon sink, *Geophysical Research Letters*, 48(20), <https://doi.org/10.1029/2021GL095325>,
1009 2021.

1010
1011 [Friedlingstein, P., O'Sullivan, M., Jones, M. W., Andrew, R. M., Bakker, D. C. E., Hauck, J.,](#)
1012 [Landschützer, P., Le Quéré, C., Luijkx, I. T., Peters, G. P., Peters, W., Pongratz, J., Schwingshackl,](#)
1013 [C., Sitch, S., Canadell, J. G., Ciais, P., Jackson, R. B., Alin, S. R., Anthoni, P., Barbero, L., Bates,](#)
1014 [N. R., Becker, M., Bellouin, N., Decharme, B., Bopp, L., Brasika, I. B. M., Cadule, P.,](#)
1015 [Chamberlain, M. A., Chandra, N., Chau, T.-T.-T., Chevallier, F., Chini, L. P., Cronin, M., Dou,](#)
1016 [X., Enyo, K., Evans, W., Falk, S., Feely, R. A., Feng, L., Ford, D. J., Gasser, T., Ghattas, J.,](#)
1017 [Gkritzalis, T., Grassi, G., Gregor, L., Gruber, N., Gürses, Ö., Harris, I., Hefner, M., Heinke, J.,](#)
1018 [Houghton, R. A., Hurtt, G. C., Iida, Y., Ilyina, T., Jacobson, A. R., Jain, A., Jarníková, T., Jersild,](#)
1019 [A., Jiang, F., Jin, Z., Joos, F., Kato, E., Keeling, R. F., Kennedy, D., Klein Goldewijk, K., Knauer,](#)
1020 [J., Korsbakken, J. I., Körtzinger, A., Lan, X., Lefèvre, N., Li, H., Liu, J., Liu, Z., Ma, L., Marland,](#)
1021 [G., Mayot, N., McGuire, P. C., McKinley, G. A., Meyer, G., Morgan, E. J., Munro, D. R., Nakaoka,](#)
1022 [S.-I., Niwa, Y., O'Brien, K. M., Olsen, A., Omar, A. M., Ono, T., Paulsen, M., Pierrot, D., Pocock,](#)
1023 [K., Poulter, B., Powis, C. M., Rehder, G., Resplandy, L., Robertson, E., Rödenbeck, C., Rosan, T.,](#)
1024 [M., Schwinger, J., Séférian, R., Smallman, T. L., Smith, S. M., Sospedra-Alfonso, R., Sun, Q.,](#)
1025 [Sutton, A. J., Sweeney, C., Takao, S., Tans, P. P., Tian, H., Tilbrook, B., Tsujino, H., Tubiello, F.,](#)
1026 [van der Werf, G. R., van Ooijen, E., Wanninkhof, R., Watanabe, M., Wimart-Rousseau, C., Yang,](#)
1027 [D., Yang, X., Yuan, W., Yue, X., Zaehle, S., Zeng, J., and Zheng, B.:](#) *Global Carbon Budget 2023,*
1028 *Earth Syst. Sci. Data*, 15, 5301–5369, <https://doi.org/10.5194/essd-15-5301-2023>, 2023.

1029 [Fyfe, J. C., Derksen, C., Mudryk, L., Flato, G. M., Santer, B. D., Swart, N. C., Molotch, N. P.,](#)
1030 [Zhang, X., Wan, H., Arora, V. K., Scinocca, J., and Jiao, Y.:](#) Large near-term projected snowpack
1031 loss over the western United States, *Nature communications*, 8(1), 14996,
1032 <https://doi.org/10.1038/ncomms14996>, 2017.

1033 Gloege, L., McKinley, G. A., Landschützer, P., Fay, A. R., Frolicher, T. L., and Fyfe, J. C.:
1034 Quantifying Errors in Observationally Based Estimates of Ocean Carbon Sink Variability, *Global*
1035 *Biogeochemical Cycles*, 35(4), <https://doi.org/10.1029/2020gb006788>, 2021.

Formatted: Font: (Default) Times New Roman, 12 pt

Deleted: Friedlingstein, P., Jones, M. W., O'Sullivan, M., Andrew, R. M., Bakker, D. C., Hauck, J., Le Quéré, C., Peters, G. P., Peters, W., Pongratz, J., Sitch, S., Canadell, J. G., Ciais, P., Jackson, R. B., Alin, S. R., Anthoni, P., Bates, N. R., Becker, M., Bellouin, N., Bopp, L., Chau, T. T. T., Chevallier, F., Chini, L. P., Cronin, M., Currie, K. I., Decharme, B., Djeutchouang, L., Dou, X., Evans, W., Feely, R. A., Feng, L., Gasser, T., Gilfillan, D., Gkritzalis, T., Grassi, G., Gregor, L., Gruber, N., Gürses, Ö., Harris, I., Houghton, R. A., Hurtt, C. C., Iida, Y., Ilyina, T., Luijkx, I. T., Jain, A. K., Jones, S. D., Kato, E., Kennedy, D., Goldewijk, K. K., Knauer, J., Korsbakken, J. A., Körtzinger, A., Landschützer, P., Lauvset, S. K., Lefèvre, N., Lienert, S., Liu, J., Marland, G., McGuire, P. C., Melton, J. R., Munro, D. R., Nabel, J. E. M. S., Nakaoka, S.-I., Niwa, Y., Ono, T., Pierrot, D., Poulter, B., Rehder, G., Resplandy, L., Robertson, E., Rödenbeck, C., Rosan, T. M., Schwinger, J., Schwingshackl, C., Séférian, R., Sutton, A. J., Sweeney, C., Tanhua, T., Tans, P. P., Tian, H., Tilbrook, B., Tubiello, F., Werf, G. V. D., Vuichard, N., Wada, C., Wanninkhof, R., Watson, A., Willis, D., Wiltshire, A. J., Yuan, W., Yue, C., Yue, X., Zaehle, S., and Zeng, J.: Global carbon budget 2021, *Earth System Science Data*, 14(4), 1917-2005, <https://doi.org/10.5194/essd-14-1917-2022>, 2022

1060 Gloege, L., Yan, M., Zheng, T. and McKinley, G. A.: Improved quantification of ocean carbon
1061 uptake by using machine learning to merge global models and pCO₂ data, *Journal of Advances in*
1062 *Modeling Earth Systems*, 14(2), <https://doi.org/10.1029/2021MS002620>, 2022.

1063

1064 Good, S. A., Martin, M., and Rayner, N. A.: EN4: Quality controlled ocean temperature and
1065 salinity profiles and monthly objective analyses with uncertainty estimates, *Journal of*
1066 *Geophysical Research Oceans*, 118(12), 6704-6717, <https://doi.org/10.1002/2013JC009067>,
1067 2013.

1068

1069 Gray, A. R., Johnson, K. S., Bushinsky, S. M., Riser, S. C., Russell, J. L., Talley, L. D.,
1070 Wanninkhof, R., Williams, N. L., and Sarmiento, J. L.: Autonomous biogeochemical floats detect
1071 significant carbon dioxide outgassing in the high-latitude Southern Ocean, *Geophysical Research*
1072 *Letters*, 45(17), 9049-9057, <https://doi.org/10.1029/2018GL078013>, 2018.

1073 Gregor, L., Lebehot, A. D., Kok, S., and Monteiro, P. M. S.: A comparative assessment of the
1074 uncertainties of global surface ocean CO₂ estimates using a machine-learning ensemble (CSIR-
1075 ML6 version 2019a) – have we hit the wall?, *Geoscientific Model Development*, 12, 5113-5136,
1076 <https://doi.org/10.5194/gmd-12-5113-2019>, 2019.

1077 Gregor, L. and Fay, A. R.: Air-sea CO₂ fluxes for surface pCO₂ data products using a standardized
1078 approach, Zenodo [code], <https://doi.org/10.5281/zenodo.5482547>, 2021.

1079 Gruber, N., Landschützer, P., and Lovenduski, N. S.: The variable Southern Ocean carbon sink,
1080 *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)
1081 [121916-063407](https://doi.org/10.1146/annurev-marine-121916-063407), 2019.

1082 [Hauck, J., Nissen, C., Landschützer, P., Rödenbeck, C., Bushinsky, S., and Olsen, A.: Sparse](#)
1083 [observations induce large biases in estimates of the global ocean CO₂ sink: and ocean model](#)
1084 [subsampling experiment, *Philosophical Transactions Of the Royal Society A*, 381:20220063,](#)
1085 <https://doi.org/10.1098/rsta.2022.0063>, 2023.

1086 Kay, J. E., Deser, C., Phillips, A., Mai, A., Hannay, C., Strand, G., Arblaster, J. M., Bates, S. C.,
1087 Danabasoglu, G., Edwards, J., Holland, M., Kushner, P., Lamarque, J-F., Lawrence, D., Lindsay,
1088 K., Middelton, A., Munoz, E., Neale, R., Oleson, K., Polvani, L., and Vertenstein, M.: The

Formatted: Font: 12 pt

Formatted: Font: 12 pt

1089 Community Earth System Model (CESM) large ensemble project: A community resource for
1090 studying climate change in the presence of internal climate variability, *Bulletin of the American*
1091 *Meteorological Society*, 96(8), 1333-1349, <https://doi.org/10.1175/BAMS-D-13-00255>, 2015.

1092 Khatiwala, S., Primeau, F., and Hall., T.: Reconstruction of the history of anthropogenic CO₂
1093 concentrations in the ocean, *Nature*, 462(7271), 346-349, <https://doi.org/10.1038/nature08526>,
1094 2009.

1095 Landschützer, P., Gruber, N., Bakker, D. C. E., and Schuster, U.: Recent variability of the global
1096 ocean carbon sink, *Global Biogeochemical Cycles*, 28(9), 927-949,
1097 <https://doi.org/10.1002/2014GB004853>, 2014.

1098 Landschützer, P., Gruber, N., Haumann, F. A., Rödenbeck, C., Bakker, D. C. E., Van Heuven, S.,
1099 Hoppema, M., Metzl, N., Sweeney, C., Takahashi, T., Brook, B., and Wanninkhof, R.: The
1100 reinvigorator of the Southern Ocean carbon sink, *Science*, 349(6253), 1221-1224,
1101 <https://doi.org/10.1126/science.aab2620>, 2015.

1102 Landschützer, P., Tanhua, T., Behncke, J., and Keppler, L.: Sailing through the Southern Ocean
1103 seas of air-sea CO₂ flux uncertainty, *Philosophical Transactions of the Royal Society A*, 381,
1104 <https://doi.org/10.1098/rsta.2022.0064>, 2023.

1105 Lenton, A. B., Matear, R. J., and Tilbrook, B.: Design of an observational strategy for quantifying
1106 the Southern Ocean uptake of CO₂, *Global Biogeochemical Cycles*, 20, 1-11.
1107 <https://doi.org/10.1029/2005GB002620>, 2006.

1108 Lenton, A. B., Tilbrook, B., Law, R. M., Bakker, D. C. E., Doney, S. C., Gruber, N., Ishii, M.,
1109 Hoppema, M., Lovenduski, N. S., Matear, R. J., McNeil, B. I., Metzl, N., Mikaloff Fletcher, S. E.,
1110 Monteiro, P. M. S., Rödenbeck, C., Sweeney, C., and Takahashi, T.: Sea-air CO₂ fluxes in the
1111 Southern Ocean for the period 1990-2009, *Biogeosciences*, 10, 4037-4054,
1112 <https://doi.org/10.5194/bg-10-4037-2013>, 2013.

1113 Le Quéré, C., Rödenbeck, C., Buitenhuis, E. T., Conway, T. J., Lagenfelds, R., Gomez, A.,
1114 Labuschagne C., Ramonet, M., Nakazawa, T., Metzl, N., Gillett, N., and Heimann, M.: Saturation
1115 of the Southern Ocean CO₂ sink due to recent climate change, *Science*, 316(5832), 1735-1738,
1116 <https://doi.org/10.1126/science.1136188>, 2007.

Deleted: .

1118 Long, M. C., Stephens, B. B., McKain, K., Sweeney, C., Keeling, R. F., Kort, E. A., Morgan, E.
1119 J., Bent, J. D., Chandra, N., Chevallier, F., Commane, R., Daube, B. C., Krummel, P. B., Loh, Z.,
1120 Lujikx, I. T., Munro, D., Patra, P., Peters, W., Ramonet, M., Rödenbeck, C., Stavert, A., Tans, P.,
1121 and Wofsy, S. C.: Strong Southern Ocean carbon uptake evident in airborne observations, *Science*,
1122 374(6572), 1275-1280, <https://doi.org/10.1126/science.abi4355>, 2021.

1123 Mackay, N., and Watson, A.: Winter air-sea CO₂ fluxes constructed from summer observations of
1124 the polar Southern Ocean suggest weak outgassing, *Journal of Geophysical Research: Oceans*,
1125 126(5), e2020JC016600, <https://doi.org/10.1029/2020JC016600>, 2021.

1126 Mackay, N., Watson, A., Suntharalingam, P., Chen, Z., and Rödenbeck, C.: Improved winter data
1127 coverage of the Southern Ocean CO₂ sink from extrapolation of summertime observations,
1128 *Communications Earth & Environment*, 3, 265, <https://doi.org/10.1038/s43247-022-00592-6>,
1129 2022.

1130 McKinley, G. A., Fay, A. R., Eddebbbar, Y. A., Gloege, L., and Lovenduski, N. S.: External forcing
1131 explains recent decadal variability of the ocean carbon sink, *AGU Advances*, 1(2),
1132 e2019AV000149, <https://doi.org/10.1029/2019AV000149>, 2020.

1133 Mongwe, N. P., Vichi, M., and Monteiro, P. M. S.: The seasonal cycle of *p*CO₂ and CO₂ fluxes in
1134 the Southern Ocean: diagnosing anomalies in CMIP5 Earth system models, *Biogeosciences*, 15(9),
1135 2851-2872, <https://doi.org/10.5194/bg-15-2851-2018>, 2018.

1136 Monteiro, P. M. S., Gregor, L., Lévy, M., Maenner, S., Sabine, C. L., and Swart, S.: Intraseasonal
1137 variability linked to sampling alias in air-sea CO₂ fluxes in the Southern Ocean, *Geophysical*
1138 *Research Letters*, 42(20), 8507-8514, <https://doi.org/10.1002/2015GL066009>, 2015.

1139 Rodgers, K. B., Lin, J., and Frölicher, T. L.: Emergence of multiple ocean ecosystem drivers in a
1140 large ensemble suite with an Earth system model, *Biogeosciences*, 12(11), 3301-3320.
1141 <https://doi.org/10.5194/bg-12-3301-2015>, 2015.

1142 Rödenbeck, C., Bakker, D. C. E., Gruber, N., Iida, Y., Jacobson, A. R., Jones, S., Landschützer,
1143 P., Metzl, N., Nakaoka, S., Olsen, A., Park, G.-H., Peylin, P., Rodgers, K. B., Sasse T. P., Schuster,
1144 U., Shutler, J. D., Valsala, V., Wanninkhof, R., and Zeng, J.: Data-based estimates of the ocean

1145 carbon sink variability – first results of the Surface Ocean pCO₂ Mapping intercomparison
1146 (SOCOM), *Biogeosciences*, 12, 7251-7278, <https://doi.org/10.5194/bg-12-7251-2015>, 2015.

Deleted: .

1147 Sabine, C., Sutton, A., McCabe, K., Lawrence-Slavas, N., Alin, S., Feely, R., Jenkins, R., Maenner,
1148 S., Meinig, C., Thomas, J., van Ooijen, E., Passmore, A., and Tilbrook, B.: Evaluation of a new
1149 carbon dioxide system for autonomous surface vehicles, *Journal of Atmospheric and Oceanic*
1150 *Technology*, 37(8), 1305-1317, <https://doi.org/10.1175/JTECH-D-20-0010.1>, 2020.

Deleted: ¶

1151 Stamell, J., Rustagi, R. R., Gloege, L., and McKinley, G. A.: Strengths and weaknesses of three
1152 Machine Learning methods for pCO₂ interpolation, *Geoscientific Model Development*
1153 *Discussions*[preprint], doi:10.5194/gmd-2020-311, 22 October 2020.

1154 Sutton, A. J., Williams, N. L., and Tilbrook, B.: Constraining Southern Ocean CO₂ flux uncertainty
1155 using uncrewed surface vehicle observations, *Geophysical Research Letters*, 48(3),
1156 e2020GL091748, <https://doi.org/10.1029/2020GL091748>, 2021.

1157 Takahashi, T., Olafsson, J., Goddard, J. G., Chipman, D. W., and Sutherland, S. C.: Seasonal
1158 variation of CO₂ and nutrients in the high-latitude surface oceans: A comparative study, *Global*
1159 *Biogeochemical Cycles*, 7(4), 843-878, <https://doi.org/10.1029/93GB02263>, 1993.

1160 Takahashi, T., Sutherland, S. C., Wanninkhof, R., Sweeney, C., Feely, R. A., Chipman, D. W.,
1161 Hales, B., Friederich, G., Chavez, F., Sabine, C., Watson, A., Bakker, D. C. E., Schuster, U., Metzl,
1162 N., Yoshikawa-Inoue, H., Ishii, M., Midorikawa, T., Nojiri, Y., Körtzinger, A., Steinhoff, T.,
1163 Hoppema, M., Olafsson, J., Arnarson, T. S., Tilbrook, B., Johannessen, T., Olsen, A., Bellerby,
1164 R., Wong, C. S., Delille, B., Bates, N. R., and de Baar, H. J. W.: Climatological mean and decadal
1165 change in surface ocean pCO₂, and net sea-air CO₂ flux over the global oceans, *Deep Sea Research*
1166 *Part II: Topical Studies in Oceanography*, 56(8-10), 554-557,
1167 <https://doi.org/10.1016/j.dsr2.2008.12.009>, 2009.

1168 Toms, B. A., Barnes, E. A., and Ebert-Uphoff, I.: Physically interpretable neural networks for the
1169 geosciences: Applications to earth system variability, *Journal of Advances in Modeling Earth*
1170 *Systems*, 12(9), e2019MS002002, <https://doi.org/10.1029/2019MS002002>, 2020.

1171 Williams, N. L., Juranek, L. W., Feely, R. A., Johnson, K. S., Sarmiento, J. L., Talley, L. D.,
1172 Dickson, A. G., Gray, A. R., Wanninkhof, R., Russell, J. L., Riser, S. C., and Takeshita, Y.:

1175 Calculating surface ocean pCO₂ from biogeochemical Argo floats equipped with pH: An
1176 uncertainty analysis, *Global Biogeochemical Cycles*, 31(3), 591-604,
1177 <https://doi.org/10.1002/2016GB005541>, 2017.

1178 Wu, Y., Bakker, D. C. E., Achterberg, E. P., Silva, A. N., Pickup D. P., Li, X., Hartman, S.,
1179 Stappard, D., Qi, D., and Tyrrell, T.: Integrated analysis of carbon dioxide and oxygen
1180 concentrations as a quality control of ocean float data, *Communications Earth & Environment*, 3,
1181 92, <https://doi.org/10.1038/s43247-022-00421-w>, 2022.

1182

1183

1184

1185

1186

1187

Page 18: [1] Deleted Thea Hatlen Heimdal 12/14/23 4:52:00 PM

▼ ◀
▲

Page 18: [1] Deleted Thea Hatlen Heimdal 12/14/23 4:52:00 PM

▼ ◀
▲

Page 18: [1] Deleted Thea Hatlen Heimdal 12/14/23 4:52:00 PM

▼ ◀
▲

Page 18: [1] Deleted Thea Hatlen Heimdal 12/14/23 4:52:00 PM

▼ ◀
▲

Page 18: [1] Deleted Thea Hatlen Heimdal 12/14/23 4:52:00 PM

▼ ◀
▲

Page 18: [1] Deleted Thea Hatlen Heimdal 12/14/23 4:52:00 PM

▼ ◀
▲

Page 18: [1] Deleted Thea Hatlen Heimdal 12/14/23 4:52:00 PM

▼ ◀
▲

Page 18: [1] Deleted Thea Hatlen Heimdal 12/14/23 4:52:00 PM

▼ ◀
▲

Page 18: [1] Deleted Thea Hatlen Heimdal 12/14/23 4:52:00 PM

▼ ◀
▲

Page 18: [1] Deleted Thea Hatlen Heimdal 12/14/23 4:52:00 PM

▼ ◀
▲

Page 18: [1] Deleted Thea Hatlen Heimdal 12/14/23 4:52:00 PM

▼ ◀
▲

Page 22: [2] Deleted Thea Hatlen Heimdal 12/13/23 10:24:00 PM

▼ ◀

▲.....
Page 22: [2] Deleted Thea Hatlen Heimdal 12/13/23 10:24:00 PM

▼.....
▲.....
Page 22: [2] Deleted Thea Hatlen Heimdal 12/13/23 10:24:00 PM

▼.....
▲.....
Page 22: [2] Deleted Thea Hatlen Heimdal 12/13/23 10:24:00 PM

▼.....
▲.....
Page 22: [2] Deleted Thea Hatlen Heimdal 12/13/23 10:24:00 PM

▼.....
▲.....
Page 22: [2] Deleted Thea Hatlen Heimdal 12/13/23 10:24:00 PM

▼.....
▲.....
Page 22: [2] Deleted Thea Hatlen Heimdal 12/13/23 10:24:00 PM

▼.....
▲.....
Page 22: [2] Deleted Thea Hatlen Heimdal 12/13/23 10:24:00 PM

▼.....
▲.....
Page 22: [2] Deleted Thea Hatlen Heimdal 12/13/23 10:24:00 PM

▼.....
▲.....
Page 22: [2] Deleted Thea Hatlen Heimdal 12/13/23 10:24:00 PM

▼.....
▲.....