



# Global Biogeochemical Cycles



# **RESEARCH ARTICLE**

10.1029/2020GB006788

# **Key Points:**

- The Large Ensemble Testbed (LET) is a powerful tool to comprehensively assess reconstruction techniques
- The amplitude of decadal variability is overestimated by 21% globally and 39% in the Southern Ocean
- Machine learning, when supplied with sufficient data, can skillfully reconstruct ocean properties

#### **Supporting Information:**

Supporting Information may be found in the online version of this article.

#### Correspondence to:

L. Gloege, ljg2157@columbia.edu

#### Citation:

Gloege, L., McKinley, G. A., Landschützer, P., Fay, A. R., Frölicher, T. L., Fyfe, J. C., et al. (2021). Quantifying errors in observationally based estimates of ocean carbon sink variability. *Global Biogeochemical Cycles*, 35, e2020GB006788. https://doi. org/10.1029/2020GB006788

Received 14 AUG 2020 Accepted 10 FEB 2021

# **Author Contributions:**

**Conceptualization:** Lucas Gloege, Galen A. McKinley, Nicole S. Lovenduski

**Data curation:** Lucas Gloege, Galen A. McKinley

Formal analysis: Lucas Gloege Funding acquisition: Galen A. McKinley

Investigation: Lucas Gloege, Galen A. McKinley, Peter Landschützer, Amanda R. Fay, John C. Fyfe, Tatiana Ilyina, Nicole S. Lovenduski, Keith B. Rodgers, Sarah Schlunegger, Yohei Takano Methodology: Lucas Gloege, Galen A. McKinley, Peter Landschützer, Nicole S. Lovenduski

© 2021. The Authors.
This is an open access article under the terms of the Creative Commons
Attribution License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

# **Quantifying Errors in Observationally Based Estimates of Ocean Carbon Sink Variability**

Lucas Gloege<sup>1</sup>, Galen A. McKinley<sup>1</sup>, Peter Landschützer<sup>2</sup>, Amanda R. Fay<sup>1</sup>, Thomas L. Frölicher<sup>3,4</sup>, John C. Fyfe<sup>5</sup>, Tatiana Ilyina<sup>2</sup>, Steve Jones<sup>6</sup>, Nicole S. Lovenduski<sup>7</sup>, Keith B. Rodgers<sup>8,9</sup>, Sarah Schlunegger<sup>10</sup>, and Yohei Takano<sup>2,11</sup>

<sup>1</sup>Lamont-Doherty Earth Observatory, Palisades, NY, USA, <sup>2</sup>Max Planck Institute for Meteorology, Hamburg, Germany, <sup>3</sup>Climate and Environmental Physics, University of Bern, Bern, Switzerland, <sup>4</sup>Oeschger Centre for Climate Change Research, University of Bern, Bern, Switzerland, <sup>5</sup>Environment and Climate Change Canada, Victoria, BC, Canada, <sup>6</sup>University of Bergen, Bergen, Norway, <sup>7</sup>University of Colorado, Boulder, CO, USA, <sup>8</sup>Center for Climate Physics, Institute for Basic Science, Busan, South Korea, <sup>9</sup>Pusan National University, Busan, South Korea, <sup>10</sup>Princeton University, Princeton, NJ, USA, <sup>11</sup>Los Alamos National Laboratory, Los Alamos, NM, USA

**Abstract** Reducing uncertainty in the global carbon budget requires better quantification of ocean CO<sub>2</sub> uptake and its temporal variability. Several methodologies for reconstructing air-sea CO<sub>2</sub> exchange from pCO<sub>2</sub> observations indicate larger decadal variability than estimated using ocean models. We develop a new application of multiple Large Ensemble Earth system models to assess these reconstructions' ability to estimate spatiotemporal variability. With our Large Ensemble Testbed, pCO2 fields from 25 ensemble members each of four independent Earth system models are subsampled as the observations and the reconstruction is performed as it would be with real-world observations. The power of a testbed is that the perfect reconstruction is known for each of the original model fields; thus, reconstruction skill can be comprehensively assessed. We find that a neural-network approach can skillfully reconstruct air-sea CO<sub>2</sub> fluxes when it is trained with sufficient data. Flux bias is low for the global mean and Northern Hemisphere, but can be regionally high in the Southern Hemisphere. The phase and amplitude of the seasonal cycle are accurately reconstructed outside of the tropics, but longer-term variations are reconstructed with only moderate skill. For Southern Ocean decadal variability, insufficient sampling leads to a 31% (15%:58%, interquartile range) overestimation of amplitude, and phasing is only moderately correlated with known truth (r = 0.54 [0.46:0.63]). Globally, the amplitude of decadal variability is overestimated by 21% (3%:34%). Machine learning, when supplied with sufficient data, can skillfully reconstruct ocean properties. However, data sparsity remains a fundamental limitation to quantification of decadal variability in the ocean carbon sink.

# 1. Introduction

The ocean significantly modulates atmospheric  $CO_2$ , having absorbed approximately 38% of industrial-age fossil carbon emissions (Friedlingstein et al., 2020). Under high emission scenarios, the ocean sink is projected to grow and become the primary sink for anthropogenic carbon emissions over the next several centuries (Randerson et al., 2015). Under low emission scenarios, such as those that would limit global warming to  $2^{\circ}$ C, the ocean carbon sink will decline rapidly as the near-surface waters that hold the bulk of anthropogenic carbon (Gruber et al., 2019) come into equilibrium with the atmosphere (Cox, 2019; Jones et al., 2016). As the long-term response to the changing atmospheric pCO<sub>2</sub> unfolds, the ocean sink will continue to be modified on seasonal to decadal timescales by climate variability and change. Ultimately, our ability to accurately monitor the fate of anthropogenic carbon in the Earth system requires a quantification of the spatially resolved variability of the ocean carbon sink on timescales from seasonal to multidecadal. To achieve this goal, global maps of surface ocean pCO<sub>2</sub> are required, from which air-sea CO<sub>2</sub> exchange can be derived.

The direction of the air-sea  $CO_2$  flux is set by the gradient in  $pCO_2$  across the air-sea interface with additional controls from the gas transfer velocity and  $CO_2$  solubility setting the magnitude. Satellites cannot directly measure surface ocean  $pCO_2$ ; therefore, hindcast simulations with ocean models (Friedlingstein et al., 2020) and observation-based gap-filling techniques (Rödenbeck et al., 2015) are integral to providing a global

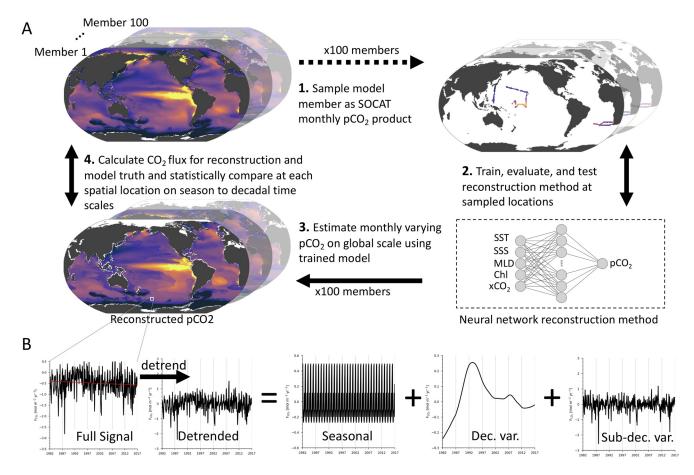
GLOEGE ET AL. 1 of 14



Project Administration: Lucas Gloege, Galen A. McKinley Software: Lucas Gloege, Peter Supervision: Galen A. McKinley Validation: Lucas Gloege Visualization: L. Gloege Writing - original draft: Lucas Gloege, Galen A. McKinley, Peter Landschützer, Amanda R. Fay, Thomas L. Frölicher, John C. Fyfe, Tatiana Ilvina, Steve Jones, Nicole S. Lovenduski, Keith B. Rodgers, Sarah Schlunegger, Yohei Takano Writing - review & editing: Lucas Gloege, Galen A. McKinley, Peter Landschützer, Amanda R. Fav. Thomas L. Frölicher, John C. Fyfe, Tatiana Ilyina, Steve Jones, Nicole S. Lovenduski, Keith B. Rodgers, Sarah Schlunegger, Yohei Takano

picture of the evolving ocean carbon sink. Essential to these techniques are high quality in-situ pCO $_2$  measurements, such as those annually compiled in the Surface Ocean CO $_2$  ATlas (SOCAT) (Bakker et al., 2016; Sabine et al., 2013). But these data are too sparse to directly constrain global air-sea CO $_2$  exchange. The latest SOCAT database release covers only 1.5% of all possible monthly 1° x 1° points from 1982 to 2019, which poses challenges to an accurate global CO $_2$  flux estimate. Current gap-filling techniques, such as the self-organizing map feed-forward neural-network (SOM-FFN)(Landschützer et al., 2016), provide continuous monthly mean estimates. However, these results lack a comprehensive, spatially resolved assessment of uncertainties. Understanding these uncertainties is important for understanding the mechanisms of variability (Landschützer et al., 2015, 2018), to compare model output to observation-based data products (Mongwe et al., 2018), to benchmark Earth system model based prediction systems (Li et al., 2019), and to assess impacts on the global carbon budget (Friedlingstein et al., 2019, 2020; Le Quéré et al., 2018). Here, we present a comprehensive, spatially resolved uncertainty assessment of the SOM-FFN method that maps pCO $_2$  from sparse observations to global coverage. Uncertainty associated with the gas transfer velocity is not accounted for in this analysis.

Our Large Ensemble Testbed uses 100 members from four Large Ensemble Earth system models, 25 members each, to evaluate the performance of the SOM-FFN over 1982–2016 given real-world pCO<sub>2</sub> sampling (Figure 1a). For each ensemble member, the pCO<sub>2</sub> reconstruction is performed in the same manner as in the SOM-FFN application to SOCAT pCO<sub>2</sub> data (see Methods). We sample the pCO<sub>2</sub> field of each testbed ensemble member as the SOCATv5 database (step 1) and use co-located driver data (see Methods) from the



**Figure 1.** The Large Ensemble Testbed. (a) Schematic of the testbed; oceanic pCO<sub>2</sub> from each of the 100 members is sampled in space and time like the SOCAT gridded product (Step 1). The sampled model output is used with auxiliary model output variables to reconstruct pCO<sub>2</sub> in the same way as the real-world application of the SOM-FFN (Step 2). pCO<sub>2</sub> is reconstructed everywhere using full-field auxiliary datasets (Step 3). Finally, CO<sub>2</sub> flux is calculated for the model truth and reconstruction for each of the 100 ensemble members and then statistically compared across seasonal to decadal time scales (Step 4). Maps in the schematic are pCO<sub>2</sub>. (b) Illustrated breakdown of CO<sub>2</sub> flux time series at a single point into seasonal, decadal, and sub-decadal variability. SOCAT, Surface Ocean CO<sub>2</sub> ATlas; SOM-FFN, self-organizing map feed-forward neural-network.

GLOEGE ET AL. 2 of 14



same ensemble member output to train, evaluate, and test the SOM-FFN (step 2). We then reconstruct full-field p $CO_2$  from the full-field driver data (step 3).  $CO_2$  flux is then calculated using the reconstructed and original Earth system model p $CO_2$  field (step 4). This is repeated for each ensemble member, providing a total of 100 unique reconstruction and model-truth pairs. To assess the performance across various timescales, we deconstruct the flux into seasonal, decadal, and sub-decadal components (Figure 1b) (see Methods). Before deconstruction, a unique point-wise linear trend is removed from each location. Performance on decadal time scales is of particular interest, since the reconstruction techniques indicate greater decadal variability than ocean models, especially in the Southern Ocean (DeVries et al., 2019; Gruber, Landschützer, & Lovenduski, 2019; Keppler & Landschützer, 2019; Ritter et al., 2017).

We emphasize that the goal of this work is not to provide an estimate of real-world air-sea  $CO_2$  exchange, but instead to assess the statistical fidelity of SOM-FFN's pCO<sub>2</sub> reconstruction, given real-world sampling. Earth system models provide plausible, though imperfect, representations of the relationships between pCO<sub>2</sub> and the associated variables required for the reconstruction. Thus, using the model output fields, we have a basis with which to test reconstruction skill on a variety of spatial and temporal scales. Fidelity is quantified by three metrics: the method's ability to capture the long-term mean, the phase, and the amplitude of seasonal to decadal time-scale variability. Our approach allows assessment of the reconstruction's fidelity across a wide range of potential states of ocean internal variability as estimated by 25 ensembles each from four independent Earth system models.

# 2. Methods

#### 2.1. SOM-FFN pCO<sub>2</sub> Interpolation

SOM-FFN (Landschützer et al., 2013, 2014, 2015) is a non-linear regression using a combination of self-organizing maps (SOM) and feed-forward neural-networks (FFN) to extrapolate from sparse pCO<sub>2</sub> observations to a global  $1^{\circ} \times 1^{\circ}$  grid at a monthly resolution. To estimate pCO<sub>2</sub> at each spatial location, SOM-FFN relies on auxiliary datasets with full, or approximately full, global coverage: Sea Surface Temperature (SST) and Surface Chlorophyll-a (Chl-a) from satellite; Sea Surface Salinity (SSS) from a compilation of *in-situ* data sources; Mixed layer depth (MLD) climatology from argo floats; and atmospheric CO<sub>2</sub> mixing ratio (xCO<sub>2</sub>). These variables serve as proxies for known processes affecting pCO<sub>2</sub>. The long-term growth of pCO<sub>2</sub> is driven by atmospheric CO<sub>2</sub> (xCO<sub>2</sub>). Solubility is set by SSS and SST. Biological uptake of dissolved inorganic carbon (DIC) is indicated by Chl-a. Biological productivity and entrainment of DIC are influenced by MLDs.

The first step uses a SOM to cluster the global ocean into 16 biogeochemical provinces based on climatological variables (surface ocean pCO<sub>2</sub> from (Takahashi et al., 2009), SST, SSS, MLD, and Chl-a). This allows for neural-network algorithms specific to each province to be developed in the second step, taking advantage of regional coherence in the dominant drivers of pCO<sub>2</sub> variability (e.g., SST in subtropics, DIC in subpolar).

The second step develops a non-linear regression to estimate pCO<sub>2</sub> given the aforementioned environmental driver variables (SST, SSS, MLD, Chl-a, and xCO<sub>2</sub>). All driver variables are monthly varying from 1982 through 2016, with the exception of climatological MLD. Any gaps in the driver data are either replaced with climatology or removed from the estimation. Within each province, a unique FFN is developed to link the driver variables to pCO<sub>2</sub> observations from SOCAT. This approach does not impose mechanistic relationships. Once the FFN algorithm is trained, tested, and evaluated on SOCAT pCO<sub>2</sub> in each province, the relationship is applied to continuous fields of driver variables to estimate pCO<sub>2</sub> at all  $1^{\circ} \times 1^{\circ}$  locations and all months from 1982 to 2016. Finally, air-sea CO<sub>2</sub> exchange is calculated following (Wanninkhof, 1992).

# 2.2. The Large Ensemble Testbed

Our 100-member Large Ensemble Testbed includes 25 randomly selected members from each of four independent initial-condition ensemble models:

- CanESM2: Second Generation Canadian Earth-System Model (RCP8.5) (Fyfe et al., 2017)
- CESM-LENS: Community Earth System Model Large Ensemble (RCP8.5) (Kay et al., 2015)

GLOEGE ET AL. 3 of 14



- GFDL-ESM2M: Geophysical Fluid Dynamics Laboratory Earth-System Model (RCP8.5) (Rodgers et al., 2015)
- MPI-GE: Max Planck Institute for Meteorology Grand Ensemble (RCP8.5) (Maher et al., 2019)

Each individual Earth system model is an imperfect representation of the actual Earth system, thus we use multiple Large Ensembles to span across the different model structures and their representation of internal variability. Each large ensemble member uses the same external forcing of historical atmospheric  $CO_2$  before 2005 and Representative Concentration Pathway 8.5 (RCP8.5) afterwards. Spread in the ensemble members is generated by perturbing the initial state of the Earth system at the start of each simulation. This is accomplished either by changing the seed value that goes into a random number generator as part of the cloud parameterization (CanESM2), perturbing the initial air-temperature field with round-off level differences (CESM-LENS), or branching off from snap-shots of the historical simulation (GFDL-ESM) or pre-industrial simulation (MPI-GE). These initial perturbations cause each ensemble member to have a unique atmosphere and ocean state at each point in time, that is, a different state of internal variability. By using many ensembles members it is possible to test the method's ability to capture the full range of  $pCO_2$  variability potential in the system under any possible climate state, not only that which occurred in the real ocean. As a specific example, the real ocean experienced an El Niño in 1997–1998. In the testbed, ensembles may have had a La Niña, El Niño or been neutral at this time. We expect only that Southern Oscillation statistics be reasonably consistent with the real world.

To create the testbed, we retrieve monthly averaged SST, SSS, Chl-a, MLD,  $xCO_2$ , and  $pCO_2$  from each member. A bilinear interpolation scheme is used to transform each field to a  $1^{\circ} \times 1^{\circ}$  rectilinear grid, the same resolution as the SOCATv5 gridded product (Sabine et al., 2013). Each member's monthly varying ocean  $pCO_2$  is then sampled at the resolution of the SOCATv5 data product, with the other variables remaining un-sampled. The sampled  $pCO_2$  field and co-located driver data for each of the 100 members constitutes the Large Ensemble Testbed capable of evaluating  $pCO_2$  interpolation methods. The intention is to create fields that mimic the environmental driver variables and SOCATv5 data used in the real-world application of the SOM-FFN interpolation. After the monthly varying  $pCO_2$  field is reconstructed for each member, air-sea  $CO_2$  exchange is calculated. The storage requirement for the 100 members testbed is  $\sim$ 500Gb, with each member, including driver data, occupying  $\sim$ 5Gb of storage. The SOM-FFN is able to reconstruct each member in about 20 minutes on a laptop. However, the computational cost will depend on the computer architecture and machine learning algorithm.

# 2.3. Air-Sea CO<sub>2</sub> Exchange

Air-sea  $CO_2$  flux is calculated in mol C m<sup>-2</sup> yr<sup>-1</sup> for each month at each 1° × 1° spatial location using the (Wanninkhof, 1992) parameterization with a scale factor of 0.27 (Sweeney et al., 2007). High-frequency output is not available for all large ensemble members. To be consistent with the flux calculation used in the real-world application of the SOM-FFN flux product, we use ERA-interim six-hourly global atmospheric reanalysis (Dee et al., 2011) as an estimate for the wind-speed variance. Saturation vapor pressure is removed from the total pressure when calculating the atmospheric partial pressure of  $CO_2$  (Dickson et al., 2007). See Text S1 for more details.

# 2.4. Temporal Decomposition

To evaluate the performance of the SOM-FFN on various time scales, an approach similar to (Cleveland et al., 1990) is used to temporally decompose the air-sea  $CO_2$  flux into additive components at each grid point (see Figure 1b for an illustration).

We first eliminate the influence of increasing atmospheric  $CO_2$  by removing a linear-trend at each  $1^{\circ} \times 1^{\circ}$  grid cell from the reconstructed air-sea  $CO_2$  flux and the model truth. Trend is calculated separately for each grid cell and is not based on the atmospheric  $CO_2$  trend. Then, a repeating seasonal cycle is calculated from the detrended time series. After removing the seasonal component, the decadal signal is isolated by applying a locally weighted regression (loess) smoother (Cleveland & Devlin, 1988) with a 10 year window. Finally, the remaining signal not explained by a linear trend, seasonal cycle, or decadal trend is here termed

GLOEGE ET AL. 4 of 14



the sub-decadal component. This decomposition was done for both the reconstructed and model truth airsea  $\rm CO_2$  flux for each of the 100 ensemble members. Statistical metrics were applied across each time scale.

#### 2.5. Statistical Metrics

The fidelity of the reconstruction is based on a suite of statistical metrics to provide a comprehensive assessment (Stow et al., 2009). Our focus is on bias, correlation, percent error in standard deviation, and average absolute error (AAE), chosen to assess if the reconstruction captures the long-term mean, temporal phasing of the signal, and variability observed in the model. Each ensemble member is treated as an equally likely climate state, and thus statistical metrics are averaged across the 100 ensemble members. Metrics are additionally calculated across the members in each Large Ensemble and the average is reported. Spread in each metric across ensemble members is quantified by the standard deviation.

Bias is calculated as the long-term mean of the reconstruction (R) minus the model truth (M),  $bias = \overline{R} - \overline{M}$ , with the overbar representing the mean over 1982–2016. Bias is a measure of the systematic discrepancy between the reconstruction and model over the long term. It is important to note that values near zero may be misleading as positive and negative discrepancies can cancel out.

Pearson correlation coefficient, r, is defined as the covariance between the reconstruction and the model divided by the product of their standard deviations,  $r = \frac{cov(R,M)}{\sigma_R \sigma_M}$ . Correlation is used to quantify the syn-

chrony between the reconstruction and model truth. Values are bounded between  $-1 \le r \le 1$ , which quantifies the degree to which reconstruction captures the phasing observed in the model. Values near 1 and -1 indicate that the reconstruction and model are perfectly in or out of phase, respectively. Intermediate values indicate a phase shift between the two signals, with values closer to zero indicating a larger phase shift between signals.

Percent error 
$$\left( \%error = \left( \frac{\sigma_R - \sigma_M}{\sigma_M} \right) * 100 \right)$$
 in the standard deviation quantifies the degree to which the

reconstruction correctly captures the amplitude of  $CO_2$  flux variability as observed in the ensemble member. This metric indicates whether the reconstruction overestimates (%error > 0), underestimates (%error < 0), or perfectly captures (%error = 0) the variability of the model truth. This metric is sensitive to the model standard deviation.

AAE quantifies how well the magnitude of variability is reconstructed in units of mol C m<sup>-1</sup> yr<sup>-1</sup>. It is defined as the absolute difference between the standard deviation of the reconstruction and of the original model field averaged across all ensemble members ( $AAE = |\sigma_R - \sigma_M|$ ).

#### 3. Results

#### 3.1. Reconstruction Bias

Regionally, the 1982-2016 mean  $CO_2$  flux from SOM-FFN can be biased high or low by more than 0.50 mol C m $^{-2}$  yr $^{-1}$  (Figure 2a), but these patches average out such that the global average bias is small (-0.01 mol Cm $^{-2}$  yr $^{-1}$ ). Regional biases are smaller in the Northern Hemisphere where data are more dense, and larger in the Indian Ocean and Southern Hemisphere where data are more sparse (Figures 2b and 2c). The mean and interquartile range of biases in the Northern and Southern Hemisphere are (0.01, -0.05:0.06) and (-0.04, -0.13:0.06) mol C m $^{-2}$  yr $^{-1}$ , respectively. Grid cells with at least 48 months of data have a mean bias that does not exceed 0.14 mol C m $^{-2}$  yr $^{-1}$  90% of the time (Figure 2c).

# 3.2. Reconstruction Phasing

Temporal correlation of the reconstruction to the original model field for each ensemble member indicates the ability of SOM-FFN to accurately capture phasing of variability at seasonal, sub-decadal, and decadal

GLOEGE ET AL. 5 of 14

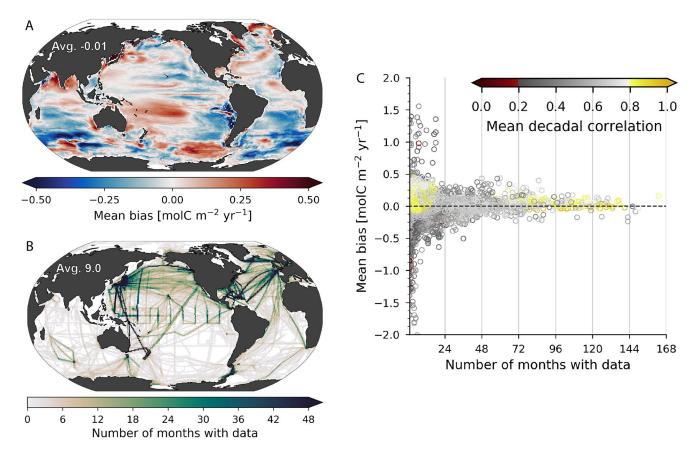


Figure 2. Reconstruction bias and sampling density. (a) Bias between reconstruction and model truth, averaged over the 100 ensemble members, each with a monthly resolution over the period 1982 through 2016. Red and blue shading indicates regions where the reconstruction is biased high or low, respectively. (b) Number of months with observations in each grid cell. The global average is displayed in each plot. (c) Cross plot of bias with number of months with data, by  $1^{\circ} \times 1^{\circ}$  grid cell. Color indicates correlation between the reconstruction and model truth on decadal time scale.

time scales (Figures 3a-3c). The standard deviation of the correlations indicates the degree to which correlations are consistent across the 100 ensemble members (Figures 3d-3f). Spatial coincidence of low standard deviations and high correlations indicates that the reconstruction performs well across all the climate states represented by the ensemble members.

Seasonally, reconstructed  $CO_2$  flux has the highest correlation to its original model field in the subtropics (Figure 3a). The large seasonal amplitude of the subtropics provides a prominent signal that the neural-network can identify (Figure S1). Higher data density in the Northern Hemisphere (Figure 2b) leads to a marginally better reconstruction which leads to better constraints on the seasonal cycle here. The lack of a prominent seasonal cycle in the tropics (Schuster et al., 2013) leads to a smaller and less coherent signal that is more difficult to reconstruct. The ability of the SOM-FFN to capture monthly variations is sporadic in the Southern Ocean and Indian Ocean, two regions that have been previously identified as having the largest mismatch toward observations and the expected seasonal amplitude increase (Landschützer et al., 2014, 2018). Despite smaller correlations around the equator and in the Southern Ocean and Indian Ocean, the global average correlation is 0.89. Additionally, regions of high correlation have low spread across the ensemble members (Figure 3d). The pattern correlation between the mean correlation (Figure 3a) and the spread of the correlations (Figure 3d) is -0.88, indicating a tight consistency between the mean result and the 100 ensemble members.

In contrast, the SOM-FFN when combined with the available observations, is less capable in reconstructing variability at sub-decadal (Figure 3b) and decadal (Figure 3c) time-scales. Global average correlation values are 0.75 and 0.58, respectively. Correlations are lower on decadal timescales (Figure 3c) than on sub-decadal timescales (Figure 3b) in the subtropics. The decadal signal is best reconstructed in the Western Pacific

GLOEGE ET AL. 6 of 14

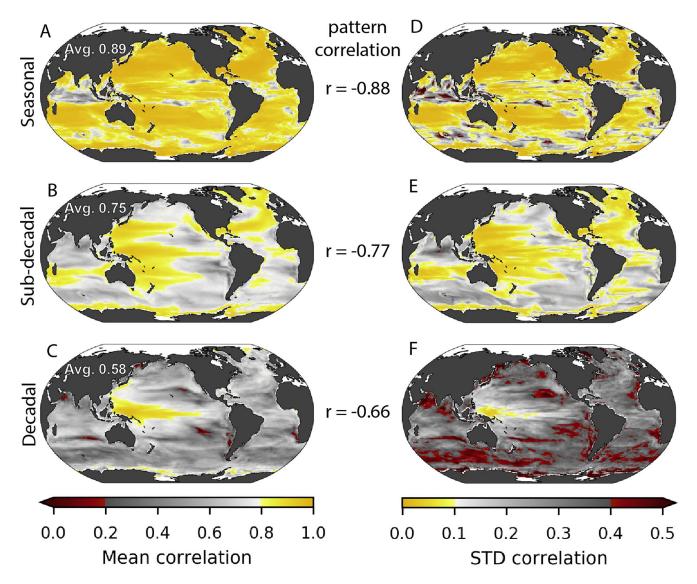


Figure 3. Phasing of SOM-FFN reconstructed variability on seasonal, sub-decadal, and decadal, compared to original model. Correlation between reconstruction and original model on (a) seasonal, (b) sub-decadal, and (c) decadal time scales, averaged across the 100 ensemble members. The global average is displayed in each plot. The standard deviation of the correlation across the 100 ensemble members is shown on (d) seasonal, (e) sub-decadal, and (f) decadal time scales. The pattern correlation between the mean and standard deviation is displayed between each pair of maps, with values close to -1 signifying high correlations are consistent across ensemble members. Note the reversed scale such that high mean correlation and low standard deviation, together indicating a robust reconstruction, have the same coloration. SOM-FFN, self-organizing map feed-forward neural-network.

warm pool. The pattern correlations between the mean and standard deviation across ensemble members are moderate (r = -0.77 for sub-decadal, and r = -0.66 for decadal), indicating a wide spread of correlations where the mean correlations are moderate. This suggests that in some ensemble members at specific locations, even the very sparse sampling that occurred was sufficient to capture the dominant modes of variation. However, this is not generally true across the ensemble, indicating a lack of robustness to the particular realization of oceanic internal variability.

# 3.3. Reconstruction of the Amplitude

Percent error of the standard deviation quantifies how well the reconstruction captures the true amplitude of variability. SOM-FFN, for the global average, overestimates the amplitude of the seasonal cycle by 7% (Figure 4a). Regionally, the reconstruction is accurate north of  $35^{\circ}N$ , but in the tropics and Southern

GLOEGE ET AL. 7 of 14

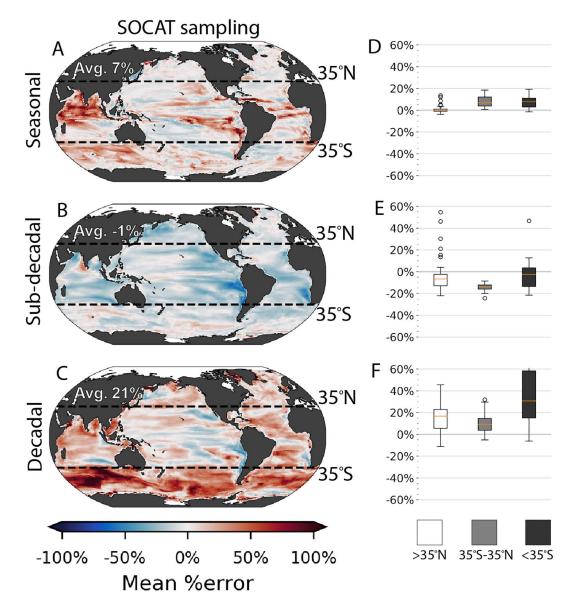


Figure 4. Error of amplitude in SOM-FFN reconstructed variability on seasonal, sub-decadal, and decadal. Percent error of  $CO_2$  flux standard deviation on (a) seasonal, (b) sub-decadal, and (c) decadal time scales, averaged across the 100 ensemble members. Global average is shown in white text. Color indicates the percentage by which the reconstruction over or under-estimates the variability. (d)–(f) Percent error as shown in (a)–(c), averaged within three regions delineated by latitude for each of the 100 ensemble members and displayed as box plots on (d) seasonal, (e) sub-decadal, and (f) decadal time scales. Boxes indicate the interquartile range (IQR), the orange line indicates the median, and circles indicate points greater than 1.5\*(IQR). IQR, interquartile range; SOM-FFN, self-organizing map feed-forward neural-network.

Hemisphere, the seasonal amplitude is overestimated by a median of 10% (3%:12%) (Figures 4a and 4d). The amplitude of sub-decadal variability is slightly underestimated at most locations, with a global average of -1% (Figure 4b).

On decadal timescales, SOM-FFN overestimates the amplitude of variability at most locations and for both the regional and global means. Globally, the overestimate is 21% (Figure 4c). In the Southern Ocean ( $<35^{\circ}$ S), the median is a 31% overestimation, with a large interquartile range across ensemble members (Figure 4f).

Percent overestimation is, by definition, inversely proportional to the model standard deviation. The fact that the four Earth system models of the Large Ensemble have different inherent amplitudes of decadal

GLOEGE ET AL. 8 of 14

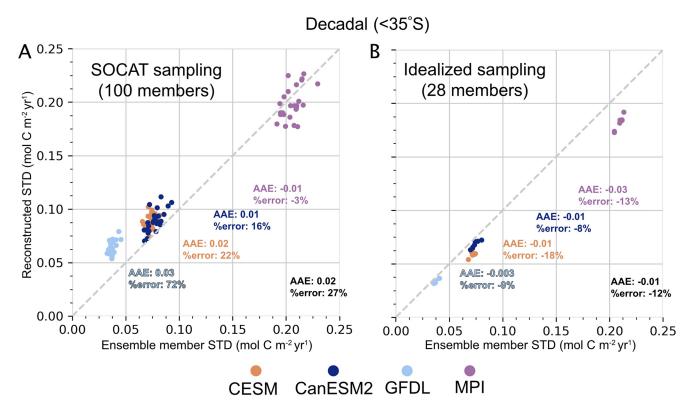


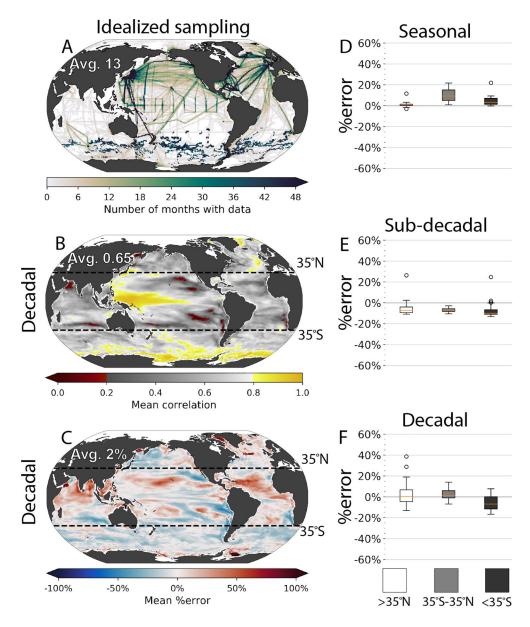
Figure 5. Cross plot of decadal standard deviation in the Southern Ocean. The reconstructed and ensemble member decadal standard deviation averaged across the Southern Ocean ( $<35^{\circ}$ S), separated by model. Colored text indicates average absolute error (AAE), and the percent error averaged across members from each model with (a) the SOCAT sampling and (b) idealized sampling. Black text indicates statistics averaged across all the ensemble members. AAE, average absolute error; SOCAT, Surface Ocean  $CO_2$  ATlas.

variability, thus influences results (Figure 5a) (Resplandy et al., 2015; Schlunegger et al., 2020). It is quite promising that the amplitude of the reconstructed decadal variability is close to its appropriate original model, as indicated by the small spread in AAE (-0.01-0.03 mol C m<sup>-2</sup> yr<sup>-1</sup>). AAE is defined as the mean of the absolute difference between the standard deviation of the reconstruction and of the original model field. Thus, the SOM-FFN is skillful in capturing the broad range of decadal variability simulated by the different Earth system models, despite the very sparse sampling. At the same time, this broad range of underlying decadal variability influences the percent error metric. MPI-GE has large decadal variability, and a low percent error (-3%); conversely, GFDL-ESM2M has small decadal variability and a high percent error (72%). Since we do not know which of these models best represent the true decadal variability of the Southern Ocean, the median across all four Large Ensembles (31%) is our best estimate for overestimation quantified as a percentage (Figure 4f). AAE of -0.01-0.03 mol C m<sup>-2</sup> yr<sup>-1</sup> is also an appropriate measure of reconstruction fidelity. However, AAE also increases as model variability declines (Figure 5a), again highlighting the challenge of not knowing which of these models best represent the true decadal variability.

# 3.4. Influence of Additional Southern Ocean Sampling

In recent years, the sampling density in the Southern Ocean has substantially increased through the launch of the fleet of drifters and Bio-Argo floats (Boutin et al., 2008; Riser et al., 2018). To assess the future impact of this new data source on our results, we test the potential impact that this additional Southern Ocean sampling would have on the reconstruction if widespread deployment had occurred for the last several decades. For this experiment, we supplement real world SOCAT sampling in the Southern Ocean (Figure 6a) for a subset of ensemble members within the Large Ensemble Testbed (Figure 1) using historical sample locations of all SOCCOM and CARIOCA measurements collapsed to a monthly climatology. These samples are assumed to have occurred at the same locations every year from 1982 to 2016. This adds 114,972 additional

GLOEGE ET AL. 9 of 14



**Figure 6.** Potential fidelity (phasing and amplitude) of SOM-FFN decadal reconstruction, had there been persistent drifters and floats in the Southern Ocean since 1982. (a) Number of months with data, with SOCAT plus idealized float sampling in the Southern Ocean; the mean (b) correlation and (c) percent error of  $CO_2$  flux standard deviation on decadal time scales across the 28 members using SOCAT plus idealized float sampling, similar to Figure 4c but with additional sampling. Box plots of percent error indicate spread among members within three regions delineated by latitude are shown on (d) seasonal, (e) sub-decadal, and (f) decadal time scales. SOCAT, Surface Ocean  $CO_2$  ATlas; SOM-FFN, self-organizing map feed-forward neural-network.

samples at 592 locations, equivalent to increasing data density globally from 1.4% with only SOCAT to 2.1% with the artificially persistent floats.

This additional sampling substantially improves the fidelity of the Southern Ocean reconstruction on all timescales (Figure 6). Enhanced sampling in the Southern Ocean also improves the reconstruction outside of the region because the biogeographic provinces of SOM-FFN are constrained by physical and biogeochemical properties, not by geography (Landschützer et al., 2014). Focusing on the Southern Ocean, the phasing of the decadal variability is improved, as indicated by higher mean correlations (Figure 6b vs. 3c). Error in the amplitude is much reduced at most locations (Figure 6c vs. 4c). The simulated additional

GLOEGE ET AL. 10 of 14



sampling also reduces the spread of amplitude error across the ensemble members on seasonal (Figure 6d), sub-decadal (Figure 6e), and decadal time scales (Figure 6f). The interquartile range for the decadal time-scale across the 28 member subset is (-11.6%, 0.0%) with a median of -6.9%. This increased sampling also substantially reduces the spread of %error estimates across the four Earth system models (Figure 5). If SO-CAT sampling had been supplemented by continuous drifters and floats in the Southern Ocean for the last 3 decades, giving only 2.1% global sampling coverage, we would now be able to reconstruct the amplitude of real-world decadal variations in the Southern Ocean carbon sink to within 20% (Figures 5b and 6f) and globally to within 2% (Figure 6c). Increasing the density of Southern Ocean observations is key to improving quantification of decadal variability in the ocean carbon sink (Bushinsky et al., 2019).

# 4. Discussion

In this work, we test the ability of the widely used SOM-FFN method to accurately reconstruct  $pCO_2$  across the global ocean. We illustrate that the reconstruction method itself can be fairly accurate across timescales, but that data sparsity remains a fundamental limitation.

These results offer the first spatially resolved quantification of the uncertainty of observation-based  $CO_2$  flux reconstruction on seasonal to decadal timescales. We address reconstruction fidelity for the ocean  $CO_2$  flux given real-world  $pCO_2$  data sparsity across a range of simulated realizations of the ocean's internal variability. We do not account for uncertainties in measurements, in the representativity of one or a small number of instantaneous  $pCO_2$  observations for a full month and a  $1^{\circ} \times 1^{\circ}$  grid cell, nor in the full-field driver data. These effects should have some impact on reconstruction fidelity, and should be assessed in future studies. Model output has previously been used to assess performance of this or similar statistical approaches for  $pCO_2$  reconstruction either using a single model (Gregor et al., 2017; Jones et al., 2015) or an ensemble of hindcast models (Lebehot et al., 2019). By incorporating many ensemble members from four Earth system models, the Large Ensemble Testbed allows for a statistically robust assessment of reconstruction performance across a range of climate states and model structures/representations. This testbed can be used to test other reconstruction approaches, as well as for development of new approaches and for evaluating new sampling strategies (Gregor et al., 2017), and is now publicly available (see Methods).

The SOM-FFN has previously been used as a reference field to assess the performance of model simulations over the historical period (Arruda et al., 2015; Bourgeois et al., 2016; Frölicher et al., 2015; Kessler & Tjiputra, 2016; Le Quéré et al., 2018; Mongwe et al., 2018). We find that SOM-FFN provides a robust global estimate of the mean  $CO_2$  uptake by the ocean, but regionally and locally, its performance is dependent on the location and the density of observations. If there are at least 48 months of data for a 35 years' timeframe, the mean bias in the long-term mean is under 0.14 mol C m<sup>-2</sup> y<sup>-1</sup> 90% of the time (Figure 2c). Mean bias can locally be much larger, particularly in poorly sampled regions such as the Southern Hemisphere. Similarly, the ability of the reconstruction to accurately capture the phase (Figure 3) and amplitude (Figure 4) of variability on sub-decadal and decadal time scales varies regionally. To improve observation-based reconstructions of the ocean carbon sink in the future, additional sampling will be critical (Figure 6).

When driven with real-world SOCAT observations and driver data, SOM-FFN indicates large amplitude decadal variability in the Southern Ocean carbon sink, with a significant slowdown in uptake over the 1990s, reaching a minimum in 2001, and then a recovery (DeVries et al., 2019; McKinley et al., 2020; Gruber, Clement, et al., 2019; Landschützer et al., 2015) until around 2011 (Keppler & Landschützer, 2019). Here, we demonstrate that, because of limited data availability, the SOM-FFN overestimates the amplitude of the decadal variability in the Southern Ocean. If the Earth system model has a larger decadal variability, the SOM-FFN provides a smaller percentage overestimate than for the Earth system models with lower amplitude decadal variability (Figure 5). However, since we do not know what the true amplitude and phasing of CO<sub>2</sub> flux variability is in the real ocean, we are unable to select one Earth system model as optimal (Mongwe et al., 2018). We find that high fidelity for seasonality (Figure 3) can co-exist with much lower fidelity for reconstruction of decadal variability (Figure 4). Our best estimate of the SOM-FFN method's overestimate of decadal variability is 31%, the median across all 100 ensemble members. A reduction of the amplitude of decadal variability by this amount would bring SOM-FFN into closer agreement with the amplitude of variability in other observation-based products (DeVries et al., 2019; Ritter et al., 2017), an ocean circula-

GLOEGE ET AL. 11 of 14



#### Acknowledgments

The authors acknowledge support from Columbia University and the Center for Climate and Life at Lamont-Doherty Earth Observatory. The authors acknowledge high-performance computing support from Cheyenne provided by NCAR's Computation and Information System Laboratory, sponsored by the National Science Foundation. L. Gloege, G. A. McKinley, A. R. Fay, and N. S. Lovenduski acknowledge funding from the National Science Foundation (OCE-1558225). P. Landschützer received funding from the European Community's Horizon 2020 project under grant agreement no. 821003 (4C). T. L. Frölicher acknowledges support from the Swiss National Science Foundation under grant PP00P2-170687 and from the European Union's Horizon 2020 Research and Innovative Program under grant agreement number 821003 (4C). Support for K. B. Rodgers was provided by the Institute for Basic Science project code IBS-R028-D1. S. Schlunegger was supported by NASA award NNX-17AI75G and NSF's Southern Ocean Carbon and Climate Observations and Modeling (SOCCOM) Project under the NSF Award PLR-1425989, with additional support from NOAA and NASA, Y. Takano was supported by the European Union's Horizon 2020 research and innovation program under grant agreement No. 641816 (CRE-SCENDO). The authors acknowledge Christian Rödenbeck for providing insightful comments during the writing of this manuscript. The Surface Ocean CO<sub>2</sub> Atlas (SOCAT) is an international effort, endorsed by the International Ocean Carbon Coordination Project (IOCCP), the Surface Ocean Lower Atmosphere Study (SOLAS), and the Integrated Marine Biosphere Research (IMBeR) program, to deliver a uniformly quality-controlled surface ocean CO<sub>2</sub> database. The many researchers and funding agencies responsible for the collection of data and quality control are thanked for their contributions to SOCAT. The authors acknowledge Data were collected and made freely available by the Southern Ocean Carbon and Climate Observations and Modeling (SOCCOM) Project funded by the National Science Foundation. Division of Polar Programs (NSF PLR -1425989), supplemented by NASA, and by the International Argo Program and the NOAA programs that contribute to it. (http://www.argo.ucsd.edu, http:// argo.jcommops.org). The Argo Program is part of the Global Ocean Observing System.

tion inverse model (DeVries et al., 2019) and hindcast ocean models (DeVries et al., 2019; Friedlingstein et al., 2019, 2020; Le Quéré et al., 2018).

Though this work indicates that SOM-FFN overestimates decadal variability of the Southern Ocean and of the globe, it does not provide a clear basis for a direct rescaling of the SOM-FFN for comparison to other estimates (Friedlingstein et al., 2020; Gruber, Clement, et al., 2019; Landschützer et al., 2015; McKinley et al., 2020). First, correlations indicate that decadal variability (Figures 3c and 3f) is only reconstructed with moderate skill in terms of phasing. Second, with respect to amplitude, the magnitude of Southern Ocean reconstructed variability from real data using SOM-FFN is 0.17 mol C m<sup>-2</sup> yr<sup>-1</sup> and from Rödenbeck et al., (2014), p. 0.16 mol C m<sup>-2</sup> yr<sup>-1</sup> (Gruber, Clement, et al., 2019). A rescaling could be derived from the mean AAE, implying a reduction of 0.02 mol C m<sup>-2</sup> yr<sup>-1</sup> to arrive at 0.14-0.15 mol C m<sup>-2</sup> yr<sup>-1</sup>. Thus, the SOM-FFN is unlikely to reconstruct a large variability if the true variability is significantly less (Figure 5). One way to constrain this range in the future could be to create a similar testbed using a suite of hindcast models, forced with realistic meteorology, that have a smaller spread in their underlying decadal variability (McKinley et al., 2020).

Great strides have been made in developing gridded synthesis products for the global carbon cycle, yet a missing component to date has been a rigorous quantitative skill assessment of these products. Using Large Ensemble output as a testbed is a powerful new approach to evaluate the skill of machine learning and other statistical extrapolations of sparse oceanographic and climatic data to global coverage fields (Deser et al., 2020). We use this approach to provide the first detailed statistical assessment of the uncertainty in a reconstruction of air-sea CO2 fluxes based on sparse in-situ ocean pCO2 data. Flux bias is low for the global mean and at most locations in the Northern Hemisphere. However, bias can be regionally high in the data-poor Southern Hemisphere. The seasonal cycle is well-captured in phase and amplitude outside of the tropics. Interannual phase and amplitude are better captured in the Northern Hemisphere and the tropics than in the Southern Hemisphere. In the Southern Ocean, insufficient sampling leads to a 31% (15%:58%) overestimation of decadal variability. Globally averaged, the amplitude of decadal variability is overestimated by 21% (3%:34%). To improve observation-based reconstructions of the ocean carbon sink, extension of sampling to include the Southern Ocean and other data-poor regions is required.

# **Conflict of Interests**

The authors declare no conflicts of interest relevant to this study.

# **Data Availability Statement**

The 100 member large ensemble testbed is publicly available at https://figshare.com/collections/Large\_ensemble\_pCO2\_testbed/4568555. Data analysis scripts are contained in GitHub repository https://github. com/lgloege/large\_ensemble\_testbed. SOCATv5 is available at https://www.socat.info/index.php/previous-versions/. ERA-interm six hourly output is available at https://apps.ecmwf.int/datasets/. Any other inquiries should be addressed to L. Gloege.

# References

Arruda, R., Calil, P. H. R., Bianchi, A. A., Doney, S. C., Gruber, N., Lima, I., & Turi, G. (2015). Air-sea CO2 fluxes and the controls on ocean surface pCO2 variability in coastal and open-ocean southwestern Atlantic Ocean: A modeling study. Biogeosciences Discussions, 12(10), 7369-7409. https://doi.org/10.5194/bgd-12-7369-2015

Bakker, D. C. E., Pfeil, B., Landa, C. S., Metzl, N., O'Brien, K. M., Olsen, A., et al. (2016). A multi-decade record of high-quality fCO2 data in version 3 of the Surface Ocean CO2 Atlas (SOCAT). Earth System Science Data, 8, 383-413. https://doi.org/10.5194/essd-8-383-2016 Bourgeois, T., Orr, J. C., Resplandy, L., Terhaar, J., Ethé, C., Gehlen, M., & Bopp, L. (2016). Coastal-ocean uptake of anthropogenic carbon. Biogeosciences, 13(14), 4167-4185. https://doi.org/10.5194/bg-13-4167-2016

Boutin, J., Merlivat, L., Hénocq, C., Martin, N., & Sallée, J. B. (2008). Air-sea CO2 flux variability in frontal regions of the Southern Ocean from carbon interface ocean atmosphere drifters. Limnology and Oceanography, 53(5part2), 2062-2079. https://doi.org/10.4319/

Bushinsky, S. M., Landschützer, P., Rödenbeck, C., Gray, A. R., Baker, D., Mazloff, M. R., et al. (2019). Reassessing Southern Ocean air-sea CO2 flux estimates with the addition of biogeochemical float observations. Global Biogeochemical Cycles, 33(11), 1370–1388. https://

Cleveland, R. B., Cleveland, W. S., McRae, J. E., & Terpenning, I. (1990). STL: A seasonal-trend decomposition. Journal of Official Statistics, 6(1), 3-73.

GLOEGE ET AL. 12 of 14



- Cleveland, W. S., & Devlin, S. J. (1988). Locally weighted regression: An approach to regression analysis by local fitting. *Journal of the American Statistical Association*, 83(403), 596–610. https://doi.org/10.1080/01621459.1988.10478639
- Cox, P. M. (2019). Emergent constraints on climate-carbon cycle feedbacks. Current Climate Change Reports, 5(4), 275–281. https://doi.org/10.1007/s40641-019-00141-y
- Dee, D. P., Uppala, S. M., Simmons, A. J., Berrisford, P., Poli, P., Kobayashi, S., et al. (2011). The ERA-interim reanalysis: Configuration and performance of the data assimilation system. *Quarterly Journal of the Royal Meteorological Society*, 137(656), 553–597. https://doi.org/10.1002/qj.828
- Deser, C., Lehner, F., Rodgers, K. B., Ault, T., Delworth, T. L., DiNezio, P. N., Fiore, A., et al. (2020). Insights from Earth system model initial-condition large ensembles and future prospects. *Nature Climate Change*, 10, 277–286. https://doi.org/10.1038/s41558-020-0731-2
- DeVries, T., Quéré, C. L., Andrews, O., Berthet, S., Hauck, J., Ilyina, T., et al. (2019). Decadal trends in the ocean carbon sink. *Proceedings of the National Academy of Sciences*, 116(24), 11646–11651. https://doi.org/10.1073/pnas.1900371116
- Dickson, A. G., Sabine, C. L., & Christian, J. R. (Eds.). (2007). Guide to best practices for ocean CO2 measurements (PICES Special Publication 3, IOCCP Report No. 8). Sidney, B. C., Canada: North Pacific Marine Science Organization.
- Friedlingstein, P., Jones, M. W., O'Sullivan, M., Andrew, R. M., Hauck, J., Peters, G. P., et al. (2019). Global Carbon Budget 2019. Earth System Science Data, 11(4), 1783–1838. https://doi.org/10.5194/essd-11-1783-2019
- Friedlingstein, P., O'Sullivan, M., Jones, M. W., Andrew, R. M., Hauck, J., Olsen, A., et al. (2020). Global Carbon Budget 2020. Earth System Science Data, 12(4), 3269–3340. https://doi.org/10.5194/essd-12-3269-2020
- Frölicher, T. L., Sarmiento, J. L., Paynter, D. J., Dunne, J. P., Krasting, J. P., & Winton, M. (2015). Dominance of the Southern Ocean in anthropogenic carbon and heat uptake in CMIP5 models. *Journal of Climate*, 28(2), 862–886. https://doi.org/10.1175/jcli-d-14-00117.1
- Fyfe, J. C., Derksen, C., Mudryk, L., Flato, G. M., Santer, B. D., Swart, N. C., et al. (2017). Large near-term projected snowpack loss over the western United States. *Nature Communications*, 8, 14996. http://doi.org/10.1038/ncomms14996
- Gregor, L., Kok, S., & Monteiro, P. M. S. (2017). Empirical methods for the estimation of Southern Ocean CO2: Support vector and random forest regression. *Biogeosciences*, 14(23), 5551–5569. https://doi.org/10.5194/bg-14-5551-2017
- Gruber, N., Clement, D., Carter, B. R., Feely, R. A., van Heuven, S., Hoppema, M., et al. (2019). The oceanic sink for anthropogenic CO2 from 1994 to 2007. Science, 363(6432), 1193–1199. https://doi.org/10.1126/science.aau5153
- Gruber, N., Landschützer, P., & Lovenduski, N. S. (2019). The variable Southern Ocean carbon sink. *Annual Review of Marine Science*, 11(1), 159–186. https://doi.org/10.1146/annurev-marine-121916-063407
- Jones, C. D., Le Quéré, C., Rödenbeck, C., Manning, A. C., & Olsen, A. (2016). Simulating the Earth system response to negative emissions. Environmental Research Letters, 11(9), 095012. https://doi.org/10.1088/1748-9326/11/9/095012
- Jones, S. D., Quéré, C. L., Rödenbeck, C., Manning, A. C., & Olsen, A. (2015). A statistical gap-filling method to interpolate global monthly surface ocean carbon dioxide data. *Journal of Advances in Modeling Earth Systems*, 7(4), 1554–1575. https://doi.org/10.1002/2014ms000416
- Kay, J. E., Deser, C., Phillips, A., Mai, A., Hannay, C., Strand, G., et al. (2015). The Community Earth System Model (CESM) large ensemble project: A community resource for studying climate change in the presence of internal climate variability. Bulletin of the American Meteorological Society, 96(8), 1333–1349. https://doi.org/10.1175/bams-d-13-00255.1
- Keppler, L., & Landschützer, P. (2019). Regional wind variability modulates the Southern Ocean carbon sink. *Scientific Reports*, 9(1), 7384. https://doi.org/10.1038/s41598-019-43826-y
- Kessler, A., & Tjiputra, J. (2016). The Southern Ocean as a constraint to reduce uncertainty in future ocean carbon sinks. *Earth System Dynamics*, 7(2), 295–312. https://doi.org/10.5194/esd-7-295-2016
- Landschützer, P., Gruber, N., & Bakker, D. C. (2016). Decadal variations and trends of the global ocean carbon sink. *Global Biogeochemical Cycles*, 30(10), 1396–1417. https://doi.org/10.1002/2015gb005359
- Landschützer, P., Gruber, N., Bakker, D. C. E., Schuster, U., Nakaoka, S., Payne, M. R., et al. (2013). A neural network-based estimate of the seasonal to inter-annual variability of the Atlantic Ocean carbon sink. *Biogeosciences*, 10(11), 7793–7815. https://doi.org/10.5194/bg-10-7793-2013
- Landschützer, P., Gruber, N., Bakker, D. C., Stemmler, I., & Six, K. D. (2018). Strengthening seasonal marine CO2 variations due to increasing atmospheric CO 2. Nature Climate Change, 8(2), 146–150. https://doi.org/10.1038/s41558-017-0057-x
- Landschützer, P., Gruber, N., Bakker, D., & Schuster, U. (2014). Recent variability of the global ocean carbon sink. *Global Biogeochemical Cycles*, 28(9), 927–949. https://doi.org/10.1002/2014gb004853
- Landschützer, P., Gruber, N., Haumann, F. A., Rödenbeck, C., Bakker, D. C. E., van Heuven, S., et al. (2015). The reinvigoration of the Southern Ocean carbon sink. *Science*, 349(6253), 1221–1224. https://doi.org/10.1126/science.aab2620
- Lebehot, A. D., Halloran, P. R., Watson, A. J., McNeall, D., Ford, D. A., Landschützer, P., et al. (2019). Reconciling observation and model trends in North Atlantic Surface CO2. Global Biogeochemical Cycles, 33(10), 1204–1222. https://doi.org/10.1029/2019gb006186
- Le Quéré, C., Andrew, R. M., Friedlingstein, P., Sitch, S., Hauck, J., Pongratz, J., et al. (2018). Global carbon budget 2018. Earth System Science Data, 10(4), 2141–2194.
- Li, H., Ilyina, T., Müller, W. A., & Landschützer, P. (2019). Predicting the variable ocean carbon sink. Science Advances, 5(4), eaav6471. https://doi.org/10.1126/sciadv.aav6471
- Maher, N., Milinski, S., Suarez-Gutierrez, L., Botzet, M., Kornblueh, L., Takano, Y., et al. (2019). The Max Planck Institute grand ensemble-enabling the exploration of climate system variability. *Journal of Advances in Modeling Earth Systems*, 11, 2050–2069. https://doi.org/10.1029/2019ms001639
- McKinley, G. A., Fay, A. R., Eddebbar, Y. A., Gloege, L., & Lovenduski, N. S. (2020). External forcing explains recent decadal variability of the ocean carbon sink. *AGU Advances*, 1(2), e2019AV000149. https://doi.org/10.1029/2019av000149
- Mongwe, N., Vichi, M., & Monteiro, P. (2018). The seasonal cycle of pCO2 and CO2 fluxes in the Southern Ocean: diagnosing anomalies in CMIP5 earth system models. *Biogeosciences*, 15(9), 2851–2872. https://doi.org/10.5194/bg-15-2851-2018
- Randerson, J. T., Lindsay, K., Munoz, E., Fu, W., Moore, J. K., Hoffman, F. M., et al. (2015). Multicentury changes in ocean and land contributions to the climate-carbon feedback. *Global Biogeochemical Cycles*, 29(6), 744–759. https://doi.org/10.1002/2014gb005079
- Resplandy, L., Séférian, R., & Bopp, L. (2015). Natural variability of CO2 and O2 fluxes: What can we learn from centuries-long climate models simulations? *Journal of Geophysical Research: Oceans*, 120(1), 384–404. https://doi.org/10.1002/2014JC010463
- Riser, S. C., Swift, D., & Drucker, R. (2018). Profiling floats in SOCCOM: Technical capabilities for studying the Southern Ocean. *Journal of Geophysical Research: Oceans*, 123(6), 4055–4073. https://doi.org/10.1002/2017JC013419
- Ritter, R., Landschützer, P., Gruber, N., Fay, A. R., Iida, Y., Jones, S., et al. (2017). Observation-based trends of the Southern Ocean carbon sink. Geophysical Research Letters, 44(24), 12339–12348. https://doi.org/10.1002/2017GL074837
- Rödenbeck, C., Bakker, D. C. E., Gruber, N., Iida, Y., Jacobson, A. R., Jones, S., et al. (2015). Data-based estimates of the ocean carbon sink variability-first results of the Surface Ocean pCO2 Mapping intercomparison (SOCOM). *Biogeosciences*, 12, 7251–7278. https://doi.org/10.5194/bg-12-7251-2015

GLOEGE ET AL. 13 of 14



- Rödenbeck, C., Bakker, D. C. E., Metzl, N., Olsen, A., Sabine, C., Cassar, N., et al. (2014). Interannual sea-air CO2 flux variability from an observation-driven ocean mixed-layer scheme. *Biogeosciences*, 11, 4599–4613. https://doi.org/10.5194/bg-11-4599-2014
- Rodgers, K. B., Lin, J., & Frölicher, T. L. (2015). Emergence of multiple ocean ecosystem drivers in a large ensemble suite with an Earth system model. *Biogeosciences*, 12(11), 3301–3320. https://doi.org/10.5194/bg-12-3301-2015
- Sabine, C. L., Hankin, S., Koyuk, H., Bakker, D. C. E., Pfeil, B., Olsen, A., et al. (2013). Surface Ocean CO2 Atlas (SOCAT) gridded data products. Earth System Science Data, 5, 145–153. https://doi.org/10.5194/essd-5-145-2013
- Schlunegger, S., Rodgers, K. B., Sarmiento, J. L., Ilyina, T., Dunne, J. P., Takano, Y., et al. (2020). Time of emergence & large ensemble intercomparison for ocean biogeochemical trends. Global Biogeochemical Cycles, 34, e2019GB006453. https://doi.org/10.1029/2019GB006453
- Schuster, U., McKinley, G. A., Bates, N., Chevallier, F., Doney, S. C., Fay, A. R., et al. (2013). An assessment of the Atlantic and Arctic sea–air CO 2 fluxes, 1990–2009. *Biogeosciences*, 10(1), 607–627. https://doi.org/10.5194/bg-10-607-2013
- Stow, C. A., Jolliff, J., McGillicuddy, D. J., Doney, S. C., Allen, J. I., Friedrichs, M. A. M., et al. (2009). Skill assessment for coupled biological/physical models of marine systems. *Journal of Marine Systems*, 76(1–2), 4–15. https://doi.org/10.1016/j.jmarsys.2008.03.011
- Sweeney, C., Gloor, E., Jacobson, A. R., Key, R. M., McKinley, G., Sarmiento, J. L., & Wanninkhof, R. (2007). Constraining global air-sea gas exchange for CO2 with recent bomb 14C measurements. *Global Biogeochemical Cycles*, 21(2), GB2015. https://doi.org/10.1029/2006gb002784
- Takahashi, T., Sutherland, S. C., Wanninkhof, R., Sweeney, C., Feely, R. A., Chipman, D. W., et al. (2009). Climatological mean and decadal change in surface ocean pCO2 and net sea–air CO2 flux over the global oceans. *Deep Sea Research Part II: Topical Studies in Oceanography*, 56(8–10), 554–577. https://doi.org/10.1016/j.dsr2.2008.12.009
- Wanninkhof, R. (1992). Relationship between wind speed and gas exchange over the ocean. *Journal of Geophysical Research: Oceans*, 97(C5), 7373–7382. https://doi.org/10.1029/92jc00188

GLOEGE ET AL. 14 of 14