

## **Supplement A: Identifying mixed-layer measurements**

Derived by wind stress and air-sea heat exchange, the mixed-layer depth (MLD) describes the maximum penetration depth of the quasi-homogeneous region of surface water [Kara et al., 2003]. Typically ranging from 20m in summer months, to 500m during the winter season in some parts of the ocean [de Boyer et al., 2004], including MLD measurements is an important additional constraint on carbon dynamics that is added from bottle measurements.

Discriminating mixed-layer measurements for each cast was conducted via a bivariate linear interpolation from a regular 2° by 2° gridded MLD climatology developed by de Boyer et al. [2004]. Their methodology was based on a change in potential density from a 10m reference measurement of 0.03 kg m<sup>-3</sup>. Approximately 900,000 CTD profiles including Argo data up to September 2008 were used to constrain their MLD climatology.

## **Supplement B: Identifying coastal data**

Carbon biogeochemical dynamics in coastal zones have been shown to be divorced from the open ocean system due to terrigenous influences [e.g., Cotrim da Cunha et al., 2007; Gibbs et al., 2006; Jickells, 1998; Seitzinger et al., 2005]. Sediment upwelling, anthropogenic influences on coastal ecosystems, and nutrient/carbon delivery from rivers have been identified as processes perturbing coastal biogeochemical dynamics from the open ocean. To mitigate these biases from our oceanic dataset, all casts with a seafloor bathymetry of 500m or less were removed from the mixed-layer training dataset. The bathymetric depth for each cast was linearly interpolated from NOAA's 1 arcminute global relief product re-gridded to 10 arcseconds [Amante and Eakins, 2009]. Eliminating coastal influences reduces the dataset by ~9%, but is important when applying the NN approach.

## **Supplement C: Anthropogenic correction for $C_T$ measurements**

The Revelle factor ( $R$ ) quantifies the relationship between the fractional change in oceanic  $p\text{CO}_2$  and  $C_T$  concentrations in an otherwise static system (Eq. S1), and is therefore a well suited empirical means to account for anthropogenic biases in  $C_T$  measurements.

$$R = \frac{C_T}{p\text{CO}_2} \frac{\Delta p\text{CO}_2}{\Delta C_T} \quad (\text{S1})$$

Rearranging equation S1 illustrates how the anthropogenic  $C_T$  component ( $\Delta C_T$ ) can be constrained if in situ  $C_T$ ,  $p\text{CO}_2$  and  $R$  are known, along with the anthropogenic change in  $p\text{CO}_2$  ( $\Delta p\text{CO}_2$ )

$$\Delta C_T = \frac{C_T}{R} \frac{\Delta p\text{CO}_2}{p\text{CO}_2} \quad (\text{S2})$$

Revelle factors and  $p\text{CO}_2$  concentrations were calculated here via the CO2SYS program developed by Pierrot et al. [2006] using bottle measurements of  $C_T$ ,  $A_T$ , temperature and salinity (phosphate and silicate concentrations were also used where available). Selection of the Mehrbach et al. [1973] constants as refitted by Dickson and Millero [1987] was based on findings by Lee et al. [2000a], McNeil et al. [2007], Millero et al. [2002], and Wanninkhof et al. [1999], and maintained consistency with the GLODAP and CARINA products [Key et al., 2004; Pierrot et al., 2010]. Here, we assume the anthropogenic rate of increase in mixed-layer  $p\text{CO}_2$  is in equilibrium with the atmosphere, which allows us to constrain  $\Delta p\text{CO}_2$  through atmospheric  $\text{CO}_2$  measurements from the Mauna Loa observation site (Dr. Pieter Tans, NOAA/ESRL, [www.esrl.noaa.gov/gmd/ccgg/trends](http://www.esrl.noaa.gov/gmd/ccgg/trends) and Dr. Ralph Keeling, Scripps Institute of Oceanography [scrippsco2.ucsd.edu/](http://scrippsco2.ucsd.edu/)). The final empirical equation to correct  $C_T$  measurements to the reference year 2000 is

$$C_{T(\text{sw},2000)} = C_{T(\text{sw},\text{in-situ year})} + \left( \frac{\text{CO}_{2(\text{atm},2000)} - \text{CO}_{2(\text{atm in-situ year})}}{p\text{CO}_{2(\text{sw},\text{in-situ year})}} \right) \frac{C_{T(\text{sw},\text{in-situ year})}}{R} \quad (\text{S3})$$

where subscripts sw and atm represent sea-water and atmosphere respectively.

Calculation of Revelle factors and  $p\text{CO}_2$  concentrations using the CO2SYS program required in situ measurements of temperature, salinity,  $A_T$  and  $C_T$ . Of the total mixed-layer  $C_T$  measurements, 8,711 (or ~28%) were missing at least one of these additional parameters required to constrain the anthropogenic correction using the proposed technique. Rather than discarding this data, 22,727 corrected  $C_T$  measurements were employed to constrain the anthropogenic correction using a 4-D linear interpolation in latitude, longitude, in situ pressure and the calculated annual anthropogenic rate of  $C_T$  increase. To evaluate the skill of the interpolation approach, we divided the 22,727 measurements into 10 equal subsets and independently interpolated the anthropogenic rate of increase. We found the approach captured the increase to within  $0.08 \mu\text{mol yr}^{-1}$  (or 8% for the mean value).

The global rate of increase in mixed-layer  $C_T$  concentration was found to be  $0.996 \mu\text{mol kg}^{-1} \text{yr}^{-1}$  (Fig. S2), which is consistent with the  $1 \mu\text{mol kg}^{-1} \text{yr}^{-1}$  anthropogenic  $C_T$  correction rate used by Lee et al. [2000b] for measurements between  $30^\circ\text{N}$  and  $30^\circ\text{S}$  and is also consistent with reported rates of increase observed at the HOT [Winn et al., 1998] and BATS [Bates, 2007] time-series stations.

Two key assumptions underlying this methodology include a constant Revelle factor over the correction period, and a global representation of atmospheric  $\text{CO}_2$  changes from the Mauna Loa site. As the ocean absorbs more anthropogenic  $\text{CO}_2$  the Revelle factor will increase, however, recent studies have estimated  $R$  to have only slightly changed over the past 2 centuries [Egleston et al., 2010], which validates our assumption of a constant  $R$  value over a maximum 20 year correction period. To evaluate the applicability of the Mauna Loa  $\Delta\text{CO}_2$  on a global scale, we compared the net change in atmospheric  $\text{CO}_2$  as observed at the Mauna Loa site to a global estimate derived from multiple stations (Thomas Conway and Pieter Tans, NOAA/ESRL, [www.esrl.noaa.gov/gmd/ccgg/trends](http://www.esrl.noaa.gov/gmd/ccgg/trends)) (Fig. S1). Here, we find a high degree of similarity between the two estimates, and when taking into consideration an uncertainty in these estimates of  $0.1 \mu\text{atm yr}^{-1}$ , the differences between the two approaches is negligible.

#### **Supplement D: Significance of anthropogenic $C_T$ correction**

To test the significance of anthropogenic  $C_T$  corrections we applied the systematic independent test approach globally (SIT, see Sect. 3) to models trained using data with and without anthropogenic  $C_T$  corrections. The global RSE for the  $C_T$  model trained using measurements without anthropogenic corrections was  $13.2 \mu\text{mol kg}^{-1}$ , or ~26% higher than the global RSE for the model trained using measurements with anthropogenic corrections ( $10.8 \mu\text{mol kg}^{-1}$ ). This difference of  $2.4 \mu\text{mol kg}^{-1}$  between the two approaches signifies the low impact of anthropogenic corrections in the models ability to constrain global  $C_T$ .

To objectively illustrate the importance of this anthropogenic correction we plot the difference between non-corrected and corrected  $C_T$  models RSE values (Eq. S4) for data in each year spanning the 30 year measurement period (Fig. S3).

$$\Delta RSE_{(yr)} = RSE_{(yr)}(\text{not corrected}) - RSE_{(yr)}(\text{corrected}) \quad (S4)$$

where yr spans the global dataset year range (i.e. 1981-2010).

The positive and increasing  $\Delta RSE_{(yr)}$  as year diverges from the reference year 2000 indicates our anthropogenic corrections enhances the global model skills. This result does not advocate that applied corrections were globally accurate, it simple confirms the importance of correcting  $C_T$  measurements to better constrain the global  $C_T$  system.

## Supplement E: Supervised SOM

A supervised form of the SOM that additionally incorporates response variable information in the clustering phase was first suggested by Kohonen [2001] and later developed by Melssen et al., [2006]. In this approach, a second neuron map of identical size to the predictor variable map established in Sect. 5.2 (wherein after referred to as the X-map) is constructed for the response variable (Y-map). Initialization of weights for the X-map remain identical to the un-supervised form, whilst the Y-map neurons are each randomly assigned a weight from within the response variable range.

Identification of the winning neuron in the X-map for data sample  $(\mathbf{x}_p, y_p)$  is determined using a distance measure that incorporates both the X-map and the Y-map

$$j(\mathbf{x}_p, y_p) = \min_j \left( (1 - \alpha(\tau)) \left[ \sum_{n=1}^N (x_{p,n} - \omega_{j,n})^2 \right]^{0.5} + \alpha(\tau) |y_p - \omega_{j(Y\text{-map})}| \right) \quad (H1)$$

where  $0 < \alpha(\tau) < 1$  is responsible for regulating the relative weight of the similarity measures of the X and Y maps. By initially setting  $\alpha(\tau)$  to 0.75, more weight is given to the neurons in the Y-map in adjusting the X-map. As  $\alpha(\tau)$  reduces linearly with iteration to 0.5, both maps are given equal weighting in identifying the winning neuron. Once the winning neuron is established, the X-map weighting vectors are updated using the same approach as presented in Sect. 5.3.

For every iteration step ( $\tau$ ), each sample is presented to the SOM model twice. In the first pass, the winning neuron in the X-map is determined and weighting vectors adjusted, whilst the second pass establishes the winning neuron in the Y-map using

$$j(\mathbf{x}_p, y_p) = \min_j \left( \alpha(\tau) \left[ \sum_{n=1}^N (x_{p,n} - \omega_{j,n})^2 \right]^{0.5} + (1 - \alpha(\tau)) |y_p - \omega_{j(\text{Y-map})}| \right) \quad (\text{H2})$$

and subsequently adjusts Y-map weighing numbers.

After the training phase is complete, response variable prediction using the X-map and any input data vector ( $\mathbf{x}_q$ ) is conducted in the same manner as presented in Sect. 5.5.

## Supplement F: Principal Component Regression

Principal Component Regression (PCR) is an empirical approach when multi-collinearity exists between predictor variables. The process (outlined in Fig. S4) first calculates the principal components ( $\mathbf{n}_1, \dots, \mathbf{n}_i, \dots, \mathbf{n}_I$ ) of the predictor variables ( $\mathbf{p}_1, \dots, \mathbf{p}_n, \dots, \mathbf{p}_N$ ). Then a least-squares multiple-linear regression is established between a subset of the principal components and the response variable ( $\mathbf{y}$ ). The subsets begin with just the first principal component, then the first two, through to all principal components. The PCR deemed optimal is simply the regression with the lowest residual standard error (RSE).

## Supplement G: Evaluating the effectiveness of a bathymetric approach in identifying coastal data

To evaluate the appropriateness of identifying coastal data based on a bathymetric depth approach, we calculated RSE values for near-coast (within 300 km of a major coastline) and open-ocean zones using the global systematic independent test predictions, however, excluding the 298 measurements already identified as terrestrially influenced and data above 70°N (Table S1). The global models ability to capture open-ocean  $A_T$  measurements was found to be ~14% (or 1.5  $\mu\text{mol kg}^{-1}$ ) better than for near-coast samples, and ~11% for  $C_T$ . This result suggests that identification of coastal measurements under a bathymetric depth approach may not be effective for ocean regions where coastal biogeochemical processes and terrestrial influences are not coupled to a shelf break, but may rather be dependent on biotic distributions. Future attempts to identify coastal measurements should therefore not solely rely on bathymetric depth.

## Supplement H: Are the neurons capturing the system?

Our optimal model configurations may be biased to the three independent datasets that constitute only 30% of the global data (see Sect. 6.1.1 Table 4). To ensure the SOM captures all important features of the global carbon system, whilst also minimising the potential influence of grouping biases, the SIT approach was applied globally using the optimal model configurations and with an increase in the SOM neuron size (Table S2).

The independent test RSE values for data below 70°N increased by 0.1 to 0.4  $\mu\text{mol kg}^{-1}$  for each step up in neuron map size (Table S2). This suggests that all important features were constrained using the three independent datasets, and that the optimal configurations remain valid on a global scale.

### **Supplement I: SOMLO model without Arctic measurements**

Uniqueness of parameter concentrations in the Arctic region (above 70°N), in particular that of salinity due to intense freshening of the water body, results in classification of Arctic measurements into features that are near exclusive to the region (Figure S5). This observation suggests Arctic measurements have little influence in constraining the remaining system.

To test this hypothesis, we compared the skill of SOMLO models trained with and without Arctic Ocean data using the SIT approach (Table S3). The skill in capturing the global carbon systems below 70°N differed by 0.1% and 2% between the two  $C_T$  and  $A_T$  models respectively, confirming that Arctic data has very little influence in the models ability to constrain the global system. This result suggests that no bias exists when comparing the skill of the global SOMLO model to the traditional MLR approach that excluded Arctic data in the regressions.

### **Supplement J: Stochastic nature of the SOM**

Initialization of neuron weights in the SOM model is a stochastic process (See Sect. 5) and can therefore lead to results that are not reproducible. In this study, the influence of this facet is dampened due to small neuron to in situ measurement ratios (1:430 for  $C_T$ ), and 800 training iteration steps converging on similar distributions of measurements among neurons for every model under static conditions.

As a test to explore stochastic influences in our model, the three independent subsets (see Sect. 6.1.1 Table 4) were each predicted 100 times using models trained under optimal configurations and the resulting RSE values examined for reproducibility (Table S4). The very small 1st standard deviation of  $0.2 \mu\text{mol kg}^{-1}$  (or 1.6%) around the mean RSE value for  $C_T$  demonstrates reproducibility in our SOMLO model and suggests a negligible influence of the stochastic SOM initialization.

## **Acknowledgements**

We would like to thank the developers of the SciPy software from which the 4D interpolations were constrained (Jones, E., Oliphant, T., Peterson, P., and others: SciPy: Open Source scientific tools for python, 2001, <http://www.scipy.org/>).

## References

Amante, C., and Eakins, B. W.: ETOPO1 1 Arc-Minute Global Relief Model: Procedures, Data Sources and Analysis, NOAA Technical Memorandum NESDIS NGDC-24, 19 pp, March, 2009.

Bates, N. R.: Interannual variability of the oceanic CO<sub>2</sub> sink in the subtropical gyre of the North Atlantic Ocean over the last 2 decades, *J. Geophys. Res.*, 112, C09013, DOI: 10.1029/2006jc003759, 2007.

Cotrim da Cunha, L., Buitenhuis, E. T., Le Quéré, C., Giraud, X., and Ludwig, W.: Potential impact of changes in river nutrient supply on global ocean biogeochemistry, *Global Biogeochem. Cycles*, 21, GB4007, DOI: 10.1029/2006gb002718, 2007.

de Boyer, M. C., Madec, G., Fischer, A. S., Lazar, A., and Iudicone, D.: Mixed layer depth over the global ocean: An examination of profile data and a profile-based climatology, *J. Geophys. Res.*, 109, C12003, DOI: 10.1029/2004jc002378, 2004.

Dickson, A. G., and Millero, F. J.: A comparison of the equilibrium constants for the dissociation of carbonic acid in seawater media, *Deep Sea Research Part A. Oceanographic Research Papers*, 34, 1733-1743, DOI: 10.1016/0198-0149(87)90021-5, 1987.

Egleston, E. S., Sabine, C. L., and Morel, F. M. M.: Revelle revisited: Buffer factors that quantify the response of ocean chemistry to changes in DIC and alkalinity, *Global Biogeochem. Cycles*, 24, GB1002, DOI: 10.1029/2008gb003407, 2010.

Gibbs, M. T., Hobday, A. J., Sanderson, B., and Hewitt, C. L.: Defining the seaward extent of New Zealand's coastal zone, *Estuarine, Coastal and Shelf Science*, 66, 240-254, DOI: 10.1016/j.ecss.2005.08.015, 2006.



Jickells, T. D.: Nutrient Biogeochemistry of the Coastal Zone, *Science*, 281, 217-222, DOI: 10.1126/science.281.5374.217 1998.

Kara, A. B., Rochford, P. A., and Hurlburt, H. E.: Mixed layer depth variability over the global ocean, *J. Geophys. Res.*, 108(C3), 3079, DOI: 10.1029/2000jc000736, 2003.

Key, R. M., Kozyr, A., Sabine, C. L., Lee, K., Wanninkhof, R., Bullister, J. L., Feely, R. A., Millero, F. J., Mordy, C., and Peng, T. H.: A global ocean carbon climatology: Results from Global Data Analysis Project (GLODAP), *Global Biogeochem. Cycles*, 18, GB4031, DOI: 10.1029/2004gb002247, 2004.

Kohonen, T.: Self-Organizing Maps, 3 ed., Number 30 in Springer Series in Information Sciences, Springer-Verlag, Berlin, 2001.

Lee, K., Millero, F. J., Byrne, R. H., Feely, R. A., and Wanninkhof, R.: The recommended dissociation constants for carbonic acid in seawater, *Geophys. Res. Lett.*, 27, 229-232, DOI: 10.1029/1999gl002345, 2000a.

Lee, K., Wanninkhof, R., Feely, R. A., Millero, F. J., and Peng, T. H.: Global relationships of total inorganic carbon with temperature and nitrate in surface seawater, *Global Biogeochem. Cycles*, 14, 979-994, DOI: 10.1029/1998GB001087, 2000b.

McNeil, B. I., Metzl, N., Key, R. M., Matear, R. J., and Corbiere, A.: An empirical estimate of the Southern Ocean air-sea CO<sub>2</sub> flux, *Global Biogeochem. Cycles*, 21, GB3011, DOI: 10.1029/2007gb002991, 2007.

Mehrbach, C., Culberson, C. H., Hawley, J. E., and Pytkowicz, R. M.: Measurement of the Apparent Dissociation Constants of Carbonic Acid in Seawater at Atmospheric Pressure, *Limnology and Oceanography*, 18, 897-907, 1973.

Melssen, W., Wehrens, R., and Buydens, L. M. C.: Supervised Kohonen networks for classification problems, *Chemometrics and Intelligent Laboratory Systems*, 83, 99-113, DOI: 10.1016/j.chemolab.2006.02.003, 2006.

Millero, F. J., Pierrot, D., Lee, K., Wanninkhof, R., Feely, R. A., Sabine, C. L., Key, R. M., and Takahashi, T.: Dissociation constants for carbonic acid determined from field measurements, *Deep Sea Research Part I: Oceanographic Research Papers*, 49, 1705-1723, DOI: 10.1016/s0967-0637(02)00093-6, 2002.

Pierrot, D., Lewis, E., and Wallace, D. W. R.: MS Excel Program Developed for CO<sub>2</sub> System Calculations, ORNL/CDIAC-105a. Carbon Dioxide Information Analysis Center, Oak Ridge National Laboratory, U.S. Department of Energy, Oak Ridge, Tennessee., DOI: 10.3334/CDIAC/otg.CO2SYS\_XLS\_CDIAC105a, 2006.

Pierrot, D., Brown, P., Van Heuven, S., Tanhua, T., Schuster, U., Wanninkhof, R., and Key, R. M.: CARINA TCO<sub>2</sub> data in the Atlantic Ocean, *Earth Syst. Sci. Data*, 2, 177-187, DOI: 10.5194/essd-2-177-2010, 2010.

Seitzinger, S. P., Harrison, J. A., Dumont, E., Beusen, A. H. W., and Bouwman, A. F.: Sources and delivery of carbon, nitrogen, and phosphorus to the coastal zone: An overview of Global Nutrient Export from Watersheds (NEWS) models and their application, *Global Biogeochem. Cycles*, 19, GB4S01, DOI: 10.1029/2005gb002606, 2005.

Wanninkhof, R., Lewis, E., Feely, R. A., and Millero, F. J.: The optimal carbonate dissociation constants for determining surface water  $p\text{CO}_2$  from alkalinity and total inorganic carbon, *Marine Chemistry*, 65, 291-301, DOI: 10.1016/s0304-4203(99)00021-3, 1999.

Winn, C. D., Li, Y.-H., Mackenzie, F. T., and Karl, D. M.: Rising surface ocean dissolved inorganic carbon at the Hawaii Ocean Time-series site, *Marine Chemistry*, 60, 33-47, DOI: 10.1016/s0304-4203(97)00085-6, 1998.

**Table S1.** Skill comparison between coastal and open-ocean predicted measurements.

| Model                | RSE <sup>a</sup> ( <i>N</i> <sup>b</sup> ) |              | % difference |
|----------------------|--|--------------|--------------|
|                      | Near-coast                                 | Open ocean   |              |
| <i>C<sub>T</sub></i> | 11.9 (4338)                                | 10.6 (18875) | 10.9         |
| <i>A<sub>T</sub></i> | 10.4 (2856)                                | 8.9 (14014)  | 14.4         |

<sup>a</sup> Residual Standard Error ( $\mu\text{mol kg}^{-1}$ )

<sup>b</sup> Number of in situ measurements

**Table S2.** RSE values ( $\mu\text{mol kg}^{-1}$ ) for models under optimal configurations and two increases in neuron map size.

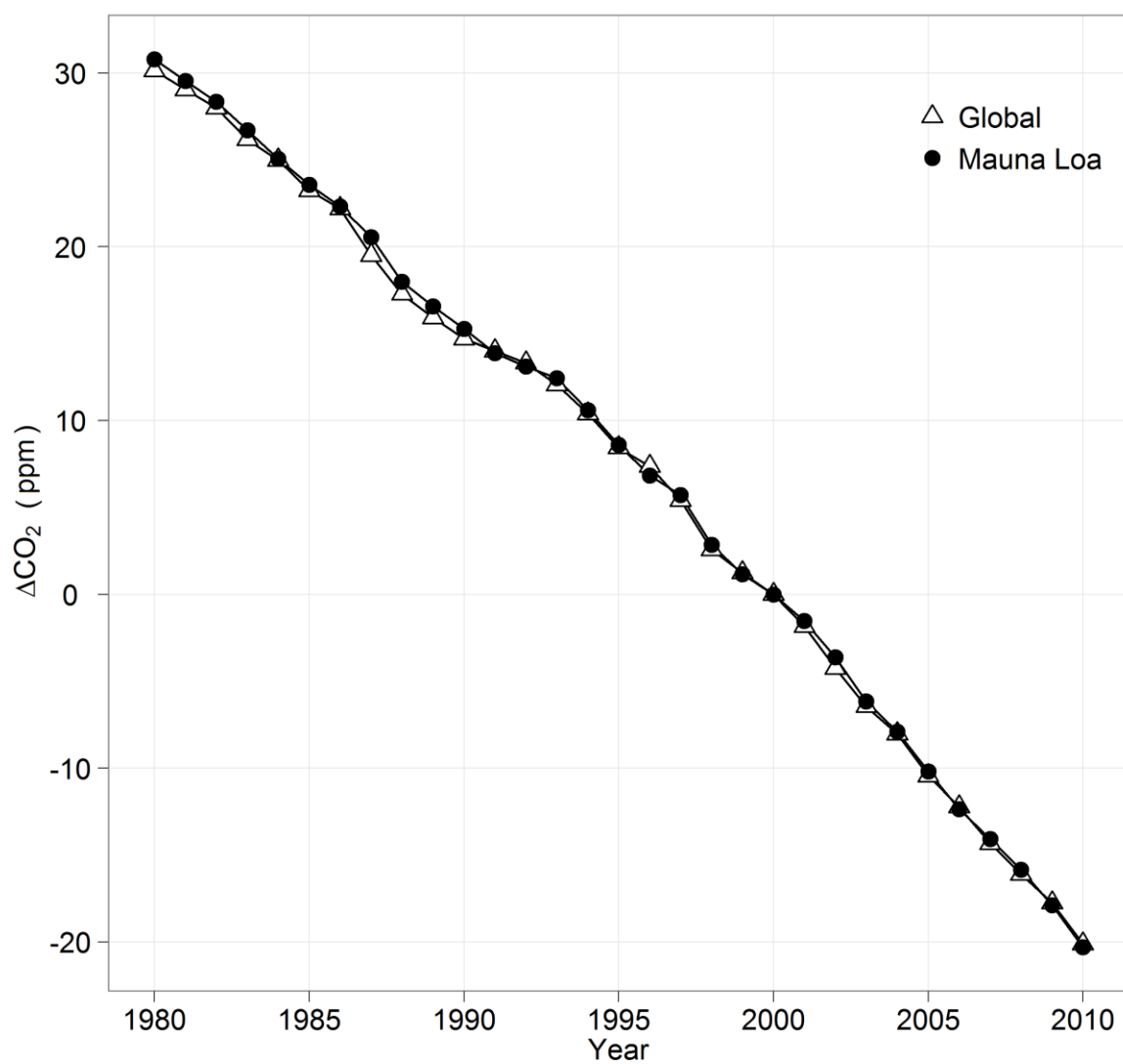
|         | $C_T$ model       |       | $A_T$ model       |       |
|---------|-------------------|-------|-------------------|-------|
|         | Number of neurons | RSE   | Number of neurons | RSE   |
| Optimal | 64                | 12.45 | 25                | 9.78  |
| Step 1  | 72                | 12.59 | 30                | 10.16 |
| Step 2  | 81                | 12.82 | 36                | 10.28 |

**Table S3.** Independent test RSE values for data below 70°N.

|       | RSE ( $\mu\text{mol kg}^{-1}$ ) |                           | % difference |
|-------|---------------------------------|---------------------------|--------------|
|       | Model with Arctic data          | Model without Arctic data |              |
| $C_T$ | 12.45                           | 12.44                     | 0.1%         |
| $A_T$ | 9.71                            | 9.9                       | 2%           |

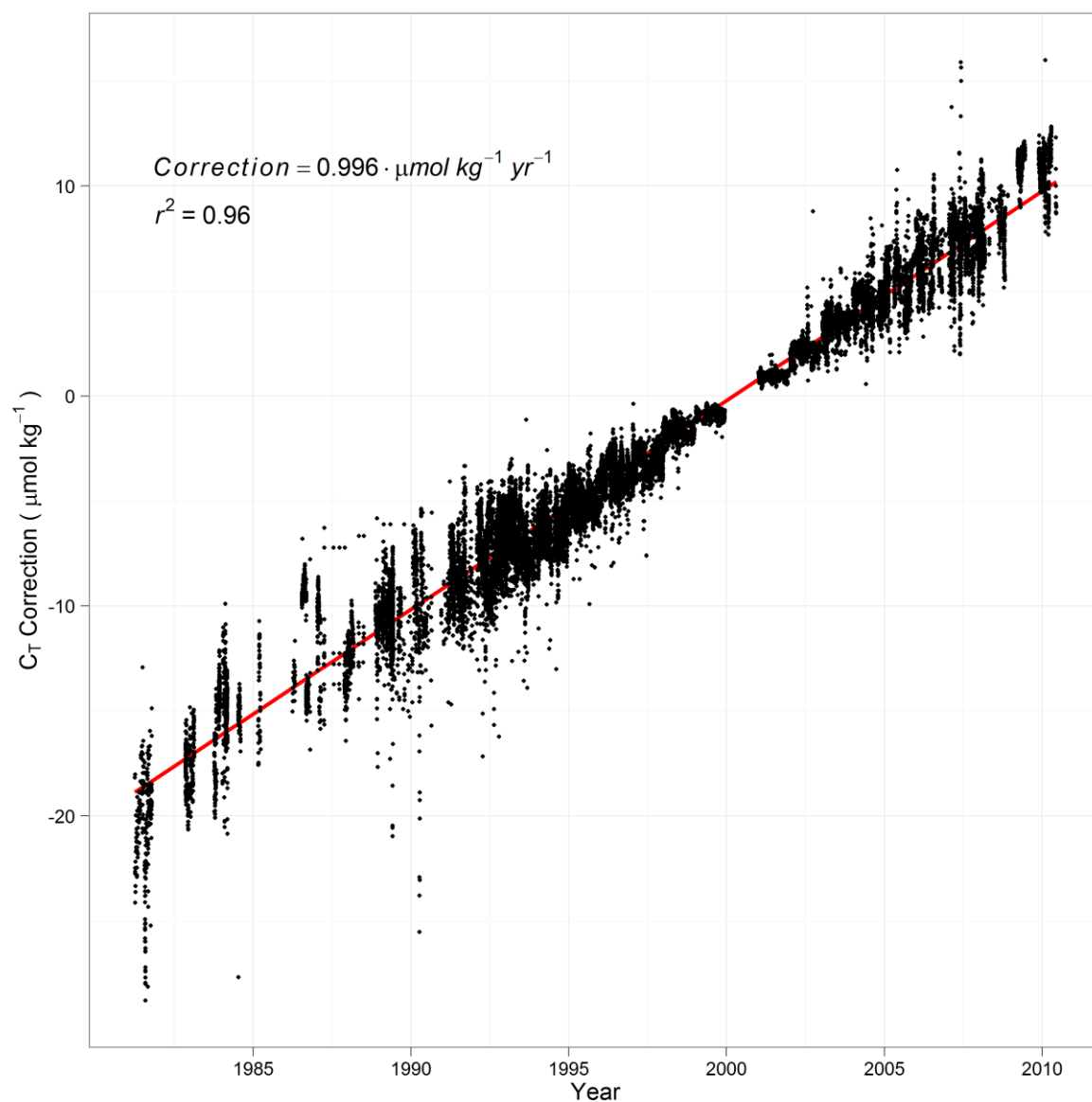
**Table S4.** RSE results for stochastic initialization test.

| Model | Mean RSE<br>( $\mu\text{mol kg}^{-1}$ ) | 1 <sup>st</sup> Standard Deviation<br>( $\mu\text{mol kg}^{-1}$ ) | % of mean |
|-------|---|---|-----------|
| $C_T$ | 12.2                                    | 0.2   | 1.6%      |
| $A_T$ | 8.2                                     | 0.1   | 1.2%      |

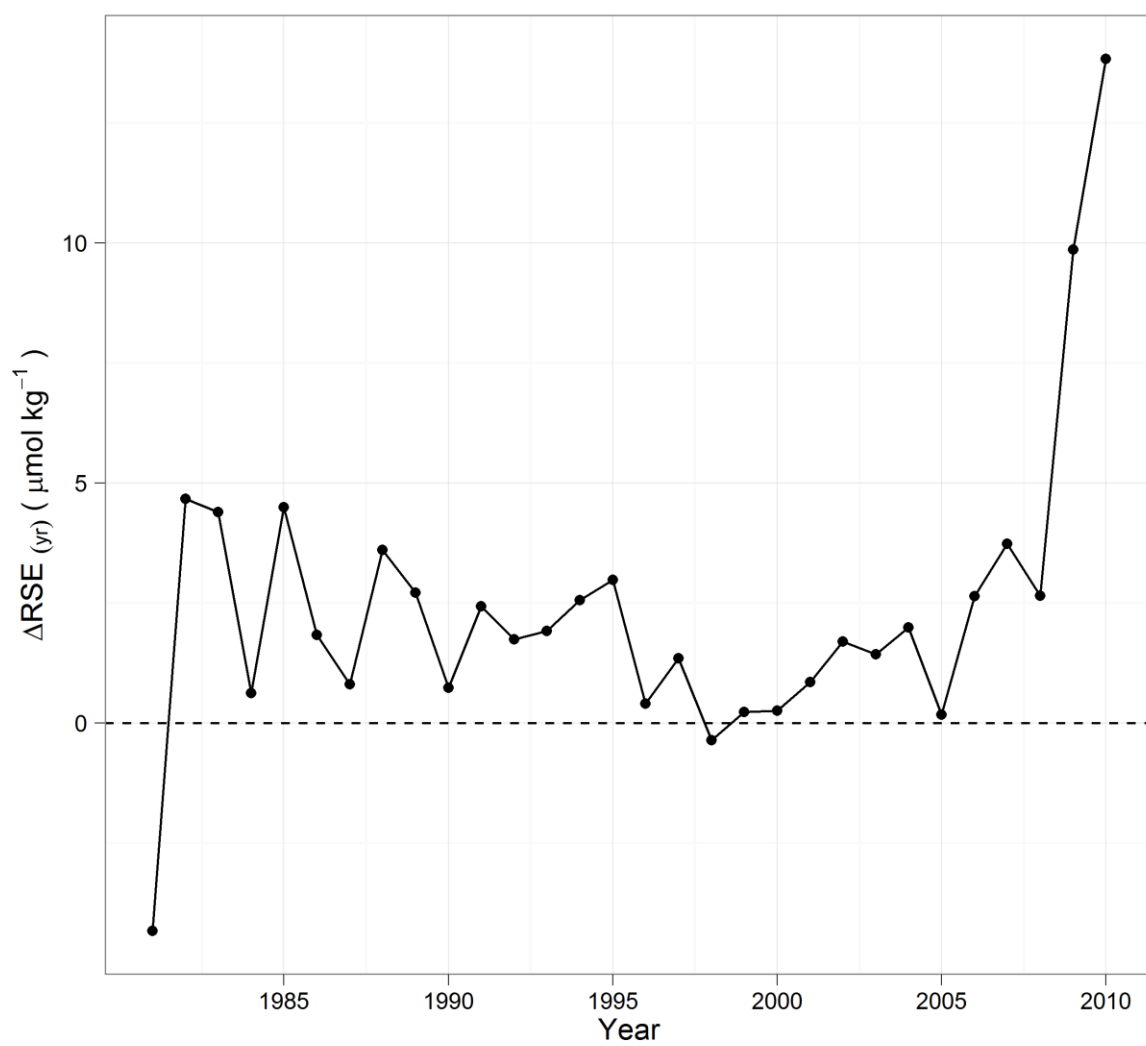


**Fig. S1.** Global and Mauna Loa site CO<sub>2</sub> difference between in situ year and the year 2000.

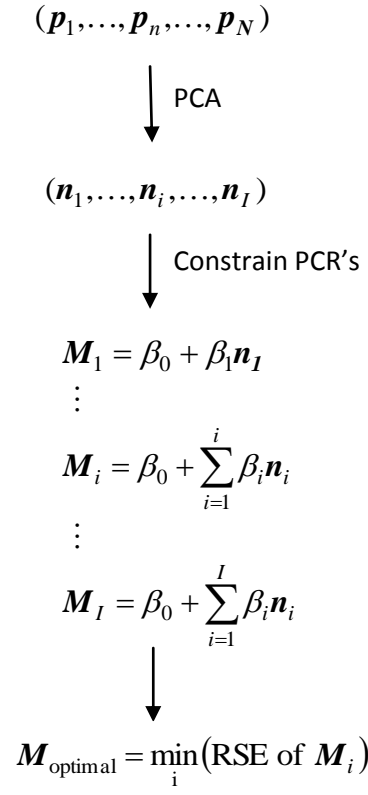




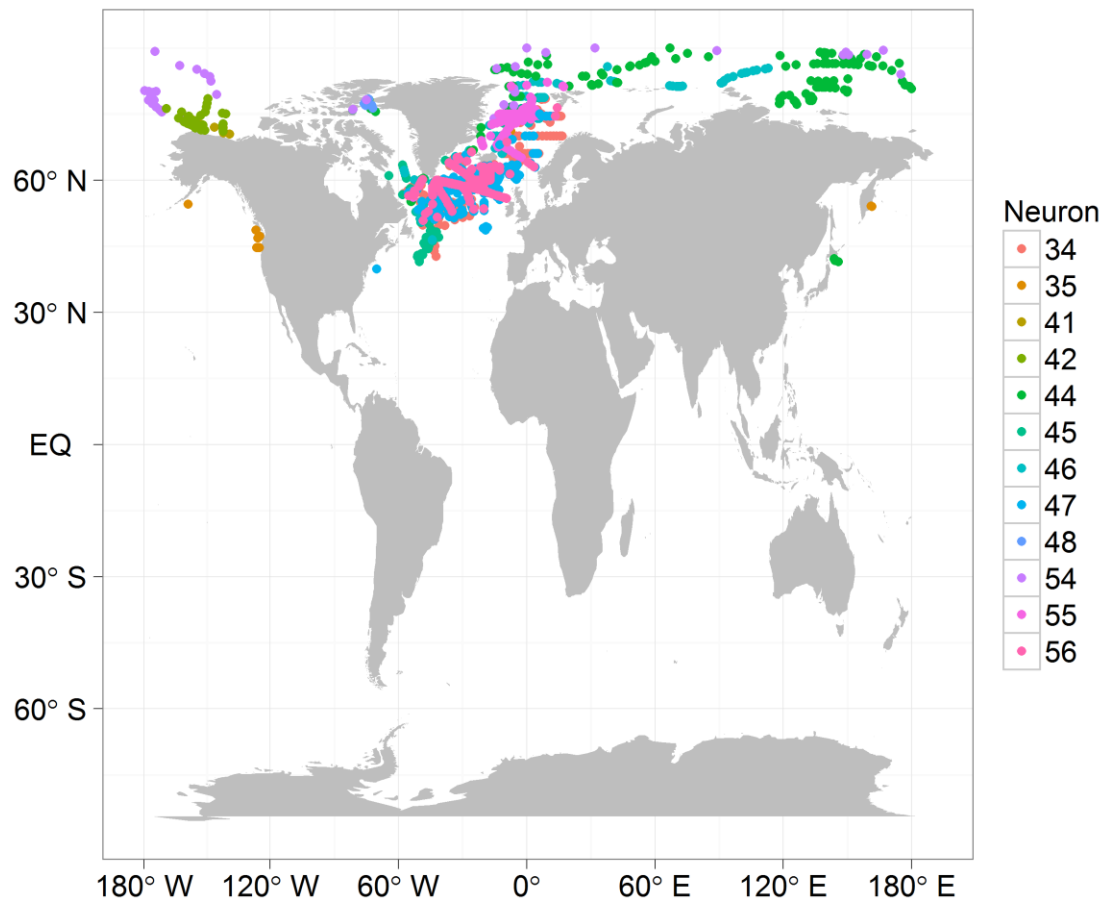
**Fig. S2.** Correction factor applied to  $C_T$  measurements defined by  $C_{T\text{correction}} = C_{T(\text{in-situ year})} - C_{T(2000)}$



**Fig. S3.** Annual  $\Delta RSE$  between  $C_T$  models trained with and without anthropogenic corrections.



**Fig. S4.** Principle Component Regression schematic.



**Fig. S5.** Distribution of measurements assigned to a neuron containing at least one Arctic measurement.