

## 1 **Supplement A: Defining the mixed-layer**

2

3 Derived by wind stress and air-sea heat exchange, the mixed-layer depth (MLD) describes the  
4 maximum penetration depth of the quasi-homogeneous region of surface water (Kara et al.,  
5 2003). Typically ranging from 20m in summer months, to 500m during the winter season in  
6 some parts of the ocean (de Boyer et al., 2004), including MLD measurements is an  
7 important additional constraint on carbon dynamics that is added from bottle measurements.

8

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

14

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

16

17 Carbon biogeochemical dynamics in coastal zones have been shown to be divorced from the  
18 open ocean system due to terrigenous influences (e.g., Cotrim da Cunha et al., 2007; Gibbs et  
19 al., 2006; Jickells, 1998; Seitzinger et al., 2005). Sediment upwelling, anthropogenic  
20 influences on coastal ecosystems, and nitrification and carbon delivery from rivers have been  
21 identified as processes perturbing coastal biogeochemical dynamics from the open ocean. To  
22 eliminate these biases from our oceanic dataset, all casts with a seafloor bathymetry of 500m  
23 or less were removed from the mixed-layer training dataset. The bathymetric depth for each  
24 cast was linearly interpolated from NOAA's 1 arcminute global relief product, re-gridded to  
25 10 arcseconds (Amante and Eakins, 2009). Eliminating coastal influences reduces the dataset  
26 by ~9%, but is important when applying the NN approach.

27

## 28 **Supplement C: Anthropogenic correction for C<sub>T</sub> measurements**

29

30 The Revelle factor (*R*) quantifies the relationship between the fractional changes of ocean  
31 *p*CO<sub>2</sub> and C<sub>T</sub> concentrations in an otherwise static system (Eq. C1), and is therefore a well  
32 suited empirical means to account for anthropogenic biases in C<sub>T</sub> measurements.

$$R = \frac{C_T}{p\text{CO}_2} \frac{\Delta p\text{CO}_2}{\Delta C_T} \quad (\text{C1})$$

Constraining the anthropogenic  $C_T$  component ( $\Delta C_T$ ) requires in-situ measurements of  $C_T$  and  $p\text{CO}_2$ , the anthropogenic change in  $p\text{CO}_2$  ( $\Delta p\text{CO}_2$ ), and Revelle factor (Eq. C2).

$$\Delta C_T = \frac{C_T}{R} \frac{\Delta p\text{CO}_2}{p\text{CO}_2} \quad (\text{C2})$$

In-situ Revelle factor and  $p\text{CO}_2$  concentrations were calculated using the CO2SYS program developed by (Pierrot et al., 2006). 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; Wanninkhof et al., 1999) and maintained consistency with 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 facilitates the ability 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

$$C_{T(\text{sw},2000)} = C_{T(\text{sw},\text{in-situ year})} + \left( \frac{\text{CO}_{2(\text{atm},2000)} - \text{CO}_{2(\text{atm},\text{in-situ year})}}{p\text{CO}_{2(\text{sw},\text{in-situ year})}} \right) \frac{C_{T(\text{sw},\text{in-situ year})}}{R} \quad (\text{C3})$$

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

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. C1). 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.

1

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

12

13 The global rate of increase in mixed-layer  $C_T$  concentration was found to be  $0.996 \mu\text{mol kg}^{-1}$   
14  $\text{yr}^{-1}$  (Fig. C2), which is consistent with the  $1 \mu\text{mol kg}^{-1} \text{yr}^{-1}$  anthropogenic  $C_T$  correction rate  
15 used by (Lee et al., 2000b) for measurements between  $30^\circ\text{N}$  and  $30^\circ\text{S}$  and is also consistent  
16 with reported rates of increase observed at the HOT (Winn et al., 1998) and BATS (Bates,  
17 2007) time-series stations.

18

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

20

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

28

29 To objectively illustrate the importance of this anthropogenic correction we plotted the  
30 difference between non-corrected and corrected  $C_T$  models RSE values (Eq. D1) for data in  
31 each year spanning the 30 year measurement period (Fig. D1).

$$32 \Delta\text{RSE}_{(\text{yr})} = \text{RSE}_{(\text{yr})}(\text{not corrected}) - \text{RSE}_{(\text{yr})}(\text{corrected}) \quad (\text{D1})$$

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

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32

The positive and increasing  $\Delta\text{RSE}_{(\text{yr})}$  as year diverges from the reference year 2000 demonstrates that the enhanced skill of the model with anthropogenic  $C_T$  corrects is a direct result of removing anthropogenic biases. This result does not advocate that applied corrections were globally accurate, it simply confirms the importance of correcting  $C_T$  measurements to better constrain the global  $C_T$  system.

### **Supplement E: Robust forward MLR based on hypothesis tests**

Following the schematic in figure E1, the routine initiates by ranking predictor variables  $(p_1, \dots, p_n, \dots, p_N)$  according to their degree of linear correlation to the response variable  $(y)$ , where  $p_{n,1}$  has the highest correlation. The primary model ( $M_1$ ) is then established by applying a least-squares multiple-linear regression (MLR) between variables  $p_{n,1}$  and  $y$  to constrain the regression coefficients  $\beta_0$  and  $\beta_{n,1}$ . In step 3, the routine expands on  $M_1$  by modelling the top two correlation ranked predictor variables ( $m = 2$ ), where  $m$  is the ranking of the modelled predictor variable with the lowest correlation to  $y$ .

To determine if multi-collinearity (MCL) exists in the expanded model ( $M_m$ ), we calculate the variance inflation factor (VIF) for each modelled variable in  $M_m$  and compare them to VIF values calculated for the same variables modelled in  $M_{m-1}$ . If the VIF value for any predictor variable  $p_{n,i}$  (where  $i < m$ ) increased by 5, it indicates existence of MLC. To reduce influences of MLC, the model is updated with interaction terms between the newly added predictor variable  $(p_{n,m})$  and any modelled variable that has a VIF increase greater than 5. To evaluate the significance of the newly added predictor variable and interaction terms, an analysis of variance (ANOVA) between the previous model ( $M_{m-1}$ ) and expanded model ( $M_m$ ) is applied. If the expanded model is found statistically better in constraining the system with a 95% confidence interval, the updates are accepted and the routine returns to step 3 where the model is again expanded with the next lowest correlation ranked predictor variable (i.e.,  $m = m + 1$ ).

If MCL is not detected, a null-hypothesis test based on the t-statistic is applied to determine if the coefficient of the new predictor variable  $(\beta_{n,m})$  is significantly different from 0 (i.e., the

1 new predictor is important in constraining the system). If it does not differ from 0 with a 95%  
2 confidence interval, the new predictor variable is defined as not significant and is  
3 subsequently rejected from the model before returning to step 3 to again expand on  $M_m$  with  
4 the next lowest ranked predictor variable.

5  
6 Once each predictor variable has had a opportunity to update the model (i.e.,  $m = I$ ), any  
7 desired higher order variable terms are incorporate into the model on the provision the first  
8 order term was found to be statistically significant. The routine then prunes the model  
9 through an iterative process removing insignificant terms based on the t-test. Once all terms  
10 are statistically significant, the final stage of the routine applies a robust MLR to the set of  
11 significant terms to reduce potential influences of outliers.

### 13 **Supplement F: Principal Component Regression**

14  
15 Principal Component Regression (PCR) is an empirical approach when multi-collinearity  
16 exists between predictor variables. The process (outlined in figure F1) first calculates the  
17 principal components  $(n_1, \dots, n_i, \dots, n_I)$  of the predictor variables  $(p_1, \dots, p_n, \dots, p_N)$ . Then  
18 a least-squares multiple-linear regression is established between a subset of the principal  
19 components and the response variable ( $y$ ). The subsets begin with just the first principal  
20 component, then the first two, through to all principal components. The PCR deemed optimal  
21 is simply the regression with the lowest residual standard error (RSE).

### 23 **Supplement G: Optimal MLR equations**

### 25 **Supplement H: Supervised SOM**

26  
27 A supervised form of the SOM that additionally incorporates response variable information in  
28 the clustering phase was first suggested by (Kohonen, 2001) and later developed by (Melssen  
29 et al., 2006). In their approach, a second neuron map of identical size to the predictor variable  
30 map established in Sect. 5.2 (wherein after referred to as the X-map) is constructed for the  
31 response variable (Y-map). Initialization of weights for the X-map remain identical to the  
32 un-supervised form, whilst the Y-map neurons are each randomly assigned a weight from  
33 within the response variable range.

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30

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

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

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 I: Geographical representation

Global position representation through latitude and longitude is problematic due to the mid-Pacific discontinuity of longitude at  $\pm 180^\circ$ , and shortening of geographical distance between degrees of longitude towards the poles. As a measure to reduce the influence of discontinuity in the model, longitude values were shifted by  $160^\circ$  West (or  $20^\circ$  East), thereby setting the  $180^\circ$  discontinuity at a position that bisects continental Africa and Europe (Fig. H1). We also tested a three dimensional n-vector calculated from latitude and longitude that eliminates both these issues.

1 The normal vector to the Earth ellipsoid (n-vector), transforms the 2-D latitude/longitude  
2 position system into a 3-D vector, whilst maintaining unique vectors for every geographical  
3 position. Employing a version of the n-vector presented by (Gade, 2010), latitude and  
4 longitude values were transformed using

$$5 \quad \mathbf{n} = \begin{bmatrix} \sin(\text{latitude}) \\ \sin(\text{longitude})\cos(\text{latitude}) \\ \cos(\text{longitude})\cos(\text{latitude}) \end{bmatrix} \quad (11)$$

6

### 7 **Supplement J: 14 region separation approach**

8

### 9 **Supplement K: Identifying poorly constrained coastal and marginal seas**

10

11 Of the 425 combined  $C_T$  and  $A_T$  predicted measurements that have a residual error greater  
12 than  $\pm 50 \mu\text{mol kg}^{-1}$ , we found that 70% (298) were located within 300 km of a major  
13 coastline. Identifying that the coastal zone around New Zealand extends up to 345 km from  
14 the shore (Gibbs et al., 2006), these anomalous independently predicted measurements are  
15 likely a result of terrestrial influences perturbing the processes affecting the carbon and  
16 hydrographic concentrations.

17

18 To evaluate the appropriateness of identifying coastal data based on a bathymetric depth  
19 approach, RSE values were calculated for coastal (within 300 km of a major coastline) and  
20 open ocean zones using the global independent test predictions, however excluding the 298  
21 measurements already identified as terrestrially influenced and data above  $70^\circ\text{N}$  (Table K1).  
22 The global models ability to capture open ocean  $A_T$  measurements is  $\sim 14\%$  (or  $1.5 \mu\text{mol kg}^{-1}$ )  
23 higher than for coastal measurements, and  $\sim 11\%$  for  $C_T$ . This result suggests that  
24 identification of coastal measurements under a bathymetric depth approach may not be  
25 effective for ocean regions where coastal biogeochemical processes and terrestrial influences  
26 are not coupled to a shelf break, but may rather be dependent on biotic distributions. Future  
27 attempts to identify coastal measurements should therefore not solely rely on bathymetric  
28 depth.

29

### 30 **Supplement L: Are the neurons capturing the system?**

31

1 Optimal model configurations may be biased to the three independent datasets that constitute  
2 only 30% of the global data (see Sect. 6.1.1 Table 4). To ensure the SOM captures all  
3 important features of the global carbon system, and minimises the potential influence of  
4 grouping biases, the GIT approach was applied globally using the optimal model  
5 configuration and with an increase in the SOM neuron size (Table L1).

6  
7 The independent test RSE values for data below 70°N increased by 0.1-0.4  $\mu\text{mol kg}^{-1}$  for  
8 each step up in neuron map size (Table L1). This suggests that all important features were  
9 constrained using the three independent datasets, and that the optimal configurations remain  
10 valid on a global scale.

### 11 12 **Supplement M: SOMLO model without Arctic measurements**

13  
14 Uniqueness of parameter concentrations in the Arctic region (above 70°N), in particular that  
15 of salinity due to intense freshening of the water body, results in classification of Arctic  
16 measurements into features that are near exclusive to the region (Figure M1). This  
17 observation suggests Arctic measurements have little influence in constraining the remaining  
18 system.

19  
20 To test this hypothesis, we compared the skill of SOMLO models trained with and without  
21 Arctic Ocean data using the GIT approach (Table M1). The skill in capturing the global  
22 carbon systems below 70°N differed by 0.1% and 2% between the two  $C_T$  and  $A_T$  models  
23 respectively, confirming that Arctic data has very little influence in the models ability to  
24 constrain the global system. This result also suggests that no bias exists in comparing the skill  
25 of the global SOMLO model to the traditional MLR approach that excluded Arctic data in the  
26 regressions.

### 27 28 **Supplement N: Stochastic nature of the SOM**

29  
30 Initialization of neuron weights in the SOM model is a stochastic process (See Sect. 5) and  
31 can therefore lead to results that are not reproducible. In this study, the influence of this facet  
32 is dampened due to small neuron to in-situ measurement ratios (1:430 for  $C_T$ ), and 800  
33 training iteration steps converging on similar distributions of measurements among neurons  
34 for every model under static conditions.



1

2 As a test to explore stochastic influences in our model, the three independent subsets (see  
3 Sect. 6.1.1 Table 4) were each predicted 100 times using models trained under optimal  
4 configurations and the resulting RSE values examined for reproducibility (Table N1). The  
5 very small 1st standard deviation of  $0.2 \mu\text{mol kg}^{-1}$  (or 1.6%) around the mean RSE value for  
6  $C_T$  demonstrates reproducibility in our SOMLO approach for constraining the carbon system  
7 and suggests a negligible influence from stochastic SOM initialization.

8

## 9 **Acknowledgements**

10

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

14

## 1 **References**

2

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

5 Bates, N. R.: Interannual variability of the oceanic CO<sub>2</sub> sink in the subtropical gyre of the North  
6 Atlantic Ocean over the last 2 decades, *J. Geophys. Res.*, 112, C09013, DOI: 10.1029/2006jc003759,  
7 2007.

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

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

14 Dickson, A. G., and Millero, F. J.: A comparison of the equilibrium constants for the dissociation of  
15 carbonic acid in seawater media, *Deep Sea Research Part A. Oceanographic Research Papers*, 34,  
16 1733-1743, DOI: 10.1016/0198-0149(87)90021-5, 1987.

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

20 Gade, K.: A Non-singular Horizontal Position Representation, *Journal of Navigation*, 63, 395 - 417,  
21 DOI: 10.1017/S0373463309990415, 2010.

22 Gibbs, M. T., Hobday, A. J., Sanderson, B., and Hewitt, C. L.: Defining the seaward extent of New  
23 Zealand's coastal zone, *Estuarine, Coastal and Shelf Science*, 66, 240-254, DOI:  
24 10.1016/j.ecss.2005.08.015, 2006.

25 Jickells, T. D.: Nutrient Biogeochemistry of the Coastal Zone, *Science*, 281, 217-222, DOI:  
26 10.1126/science.281.5374.217 1998.

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

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

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

6 Lee, K., Millero, F. J., Byrne, R. H., Feely, R. A., and Wanninkhof, R.: The recommended dissociation  
7 constants for carbonic acid in seawater, *Geophys. Res. Lett.*, 27, 229-232, DOI:  
8 10.1029/1999gl002345, 2000a.

9 Lee, K., Wanninkhof, R., Feely, R. A., Millero, F. J., and Peng, T. H.: Global relationships of total  
10 inorganic carbon with temperature and nitrate in surface seawater, *Global Biogeochem. Cycles*, 14,  
11 979–994, DOI: 10.1029/1998GB001087, 2000b.

12 McNeil, B. I., Metzl, N., Key, R. M., Matear, R. J., and Corbiere, A.: An empirical estimate of the  
13 Southern Ocean air-sea CO<sub>2</sub> flux, *Global Biogeochem. Cycles*, 21, GB3011, DOI:  
14 10.1029/2007gb002991, 2007.

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

18 Melssen, W., Wehrens, R., and Buydens, L. M. C.: Supervised Kohonen networks for classification  
19 problems, *Chemometrics and Intelligent Laboratory Systems*, 83, 99-113, DOI:  
20 10.1016/j.chemolab.2006.02.003, 2006.

21 Millero, F. J., Pierrot, D., Lee, K., Wanninkhof, R., Feely, R. A., Sabine, C. L., Key, R. M., and Takahashi,  
22 T.: Dissociation constants for carbonic acid determined from field measurements, *Deep Sea*  
23 *Research Part I: Oceanographic Research Papers*, 49, 1705-1723, DOI: 10.1016/s0967-  
24 0637(02)00093-6, 2002.

25 Pierrot, D., Lewis, E., and Wallace, D. W. R.: MS Excel Program Developed for CO<sub>2</sub> System  
26 Calculations, ORNL/CDIAC-105a. Carbon Dioxide Information Analysis Center, Oak Ridge National  
27 Laboratory, U.S. Department of Energy, Oak Ridge, Tennessee., DOI:  
28 10.3334/CDIAC/otg.CO2SYS\_XLS\_CDIA105a, 2006.

- 1 Pierrot, D., Brown, P., Van Heuven, S., Tanhua, T., Schuster, U., Wanninkhof, R., and Key, R. M.:  
2 CARINA TCO<sub>2</sub> data in the Atlantic Ocean, *Earth Syst. Sci. Data*, 2, 177-187, DOI: 10.5194/essd-2-177-  
3 2010, 2010.
- 4 Seitzinger, S. P., Harrison, J. A., Dumont, E., Beusen, A. H. W., and Bouwman, A. F.: Sources and  
5 delivery of carbon, nitrogen, and phosphorus to the coastal zone: An overview of Global Nutrient  
6 Export from Watersheds (NEWS) models and their application, *Global Biogeochem. Cycles*, 19,  
7 GB4S01, DOI: 10.1029/2005gb002606, 2005.
- 8 Wanninkhof, R., Lewis, E., Feely, R. A., and Millero, F. J.: The optimal carbonate dissociation  
9 constants for determining surface water pCO<sub>2</sub> from alkalinity and total inorganic carbon, *Marine*  
10 *Chemistry*, 65, 291-301, DOI: 10.1016/s0304-4203(99)00021-3, 1999.
- 11 Winn, C. D., Li, Y.-H., Mackenzie, F. T., and Karl, D. M.: Rising surface ocean dissolved inorganic  
12 carbon at the Hawaii Ocean Time-series site, *Marine Chemistry*, 60, 33-47, DOI: 10.1016/s0304-  
13 4203(97)00085-6, 1998.

1 **Table G1.** Ad-hoc and universal  $C_T$  regression equations with interaction terms (Int.).

	Intercept	T	S	DO	N	Si	P	Int.	Int.	Int.	Int.	2
North Pacific summer	1066.9	–	24.16	0.38	5.07	–	–	–	–	–	–	3
North Pacific winter	868.73	-7.95	36.54	–	–	4.73	–	-0.01*Si*DO	–	–	–	
Southern Ocean summer	698.74	–	35.84	0.25	–	0.42	83.28	–	–	–	–	
Southern Ocean winter	1494.14	-48.81	22.3	-0.31	–	0.48	–	0.03*T*DO	0.92*T*S	–	–	
Northwest Atlantic summer	1709.15	–	9.13	–	9.14	–	–	–	–	–	–	
Northwest Atlantic winter	1128.79	-5.93	28.71	–	17.31	–	–	-0.05*N*DO	–	–	–	
Northeast Atlantic summer	1625.17	–	12.35	–	6.13	–	–	–	–	–	–	
Northeast Atlantic winter	1013.95	61.48	32.75	–	-61.93	-2.05	-12.54	0.4*N*Si	-0.3*T*DO	-1.67*T*S	1.78*N*S	
Equatorial Pacific	467.55	-7.11	48.4	–	2.34	1.44	38.85	–	–	–	–	
Sub-tropical North Pacific summer	519.77	-9.98	48.97	–	20.25	–	1.92	–	–	–	–	
Sub-tropical North Pacific winter	236.82	-5.32	50.65	0.38	-1.5	–	139.22	–	–	–	–	
Sub-tropical South Pacific summer	147.12	-4.61	52.63	0.34	7.48	1.67	72.68	–	–	–	–	
Sub-tropical South Pacific winter	643.05	-12.01	46.68	–	–	-1.28	107.88	–	–	–	–	
Indian Ocean summer	551.82	-6.59	45.21	–	–	–	27.16	0.14*S*N	–	–	–	
Indian Ocean winter	1733.55	-1.84	–	-4.78	18.13	2.64	67.1	0.17*DO*S	–	–	–	
Sub-tropical North Atlantic summer	619.34	-8.18	46.94	-0.37	-31.07	–	–	–	–	–	–	
Sub-tropical North Atlantic winter	765	-7.36	39.89	–	–	-4.88	109.27	–	–	–	–	
Equatorial Atlantic	163.5	-8.91	59.32	-0.17	–	–	–	–	–	–	–	
Sub-tropical South Atlantic	2277.89	-6.15	–	-7.48	-5.08	–	74.92	0.2*DO*S	–	–	–	
Universal model	596.77	-8.21	45.5	-0.17	1.12	0.45	17.83	0.01*T*DO	1.52*T*P	–	–	

1 **Table G2.** Ad-hoc and universal  $A_T$  regression equations with interaction terms (Int.).

	Intercept	T	S	$S^2$	DO	Si	P	Int.	Int.	Int.
Sub-tropical Oceans	2064.66	-0.3	-47.57	1.54	0.13	-1.12	10.1	–	–	–
Equatorial Pacific	1142.6	-1.39	–	0.96	0.14	–	-3.51	–	–	–
North Atlantic	1543.52	-4.78	–	0.64	-0.04	-0.29	-9.04	0.13*S*T	–	–
North Pacific	721.6	–	44.31	–	0.09	-7.81	9.97	0.24*S*Si	–	–
Southern Ocean	7661.04	-1.46	-362.53	5.86	0.54	-12.17	-6.56	0.08*S*T	0.44*S*Si	-0.01*DO*Si
Global	1972.44	-12.78	-33.44	1.19	0.16	0.39	6.89	0.37*S*T	–	–

1 **Table K1.** Skill comparison between coastal and open ocean predicted measurements.

Model	RSE <sup>a</sup> (N <sup>b</sup> )		% difference
	Coastal	Open Ocean	
C <sub>T</sub>	11.9 (4338)	10.6 (18875)	10.9
A <sub>T</sub>	10.4 (2856)	8.9 (14014)	14.4

2

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

4 <sup>b</sup> N = number of in-situ measurements

5

- 1 **Table L1.** RSE values ( $\mu\text{mol kg}^{-1}$ ) for models under optimal configurations and two increases  
 2 in neuron map size.

	C <sub>T</sub> model		A <sub>T</sub> 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

3



1 **Table M1.** Independent test RSE values for data below 70°N.

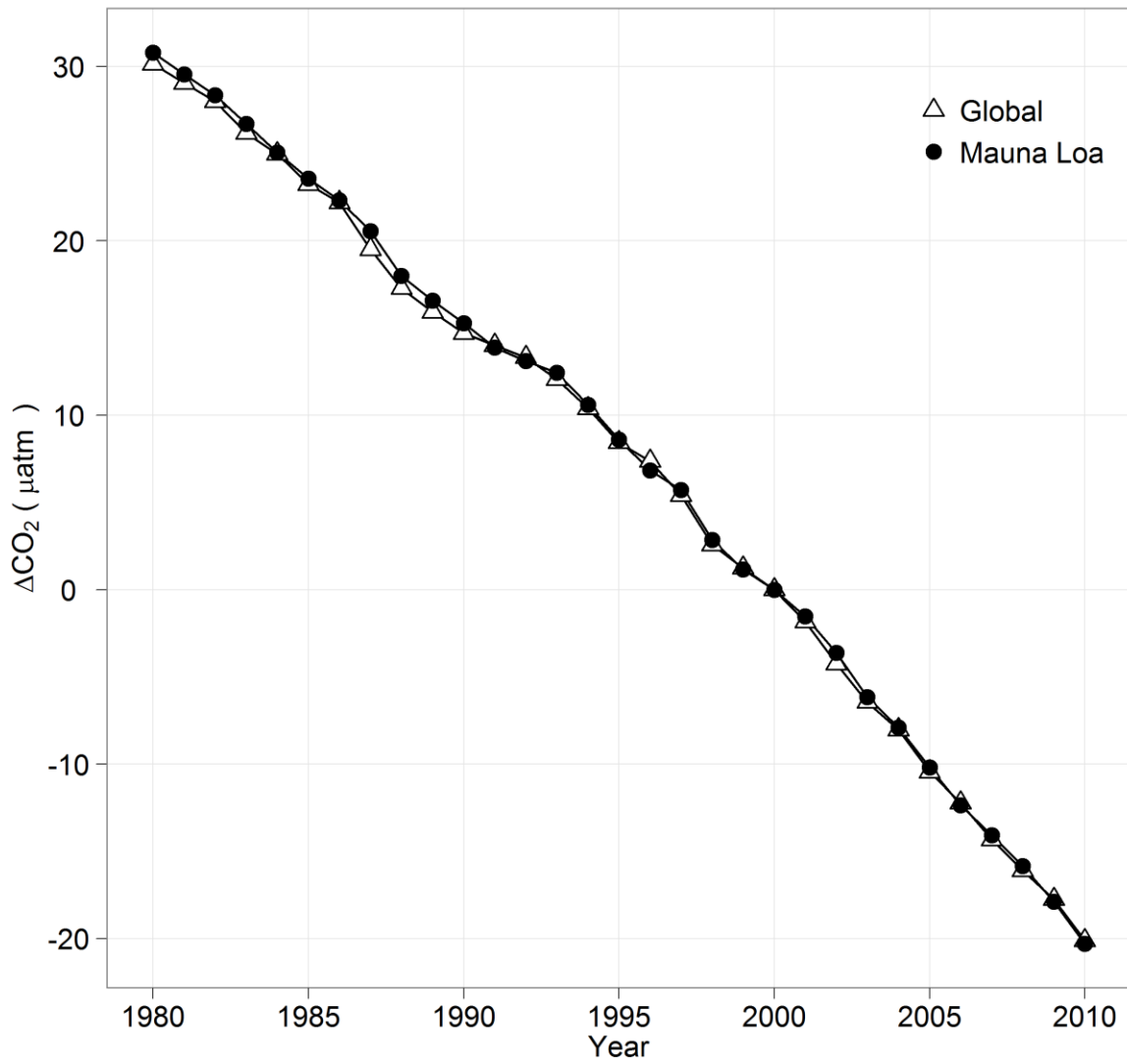
	RSE ( $\mu\text{mol kg}^{-1}$ )		% difference
	Model with Arctic data	Model without Arctic data	
C <sub>T</sub>	12.45	12.44	0.1%
A <sub>T</sub>	9.71	9.9	2%

2

1 **Table N1.** RSE results for stochastic initialization test.

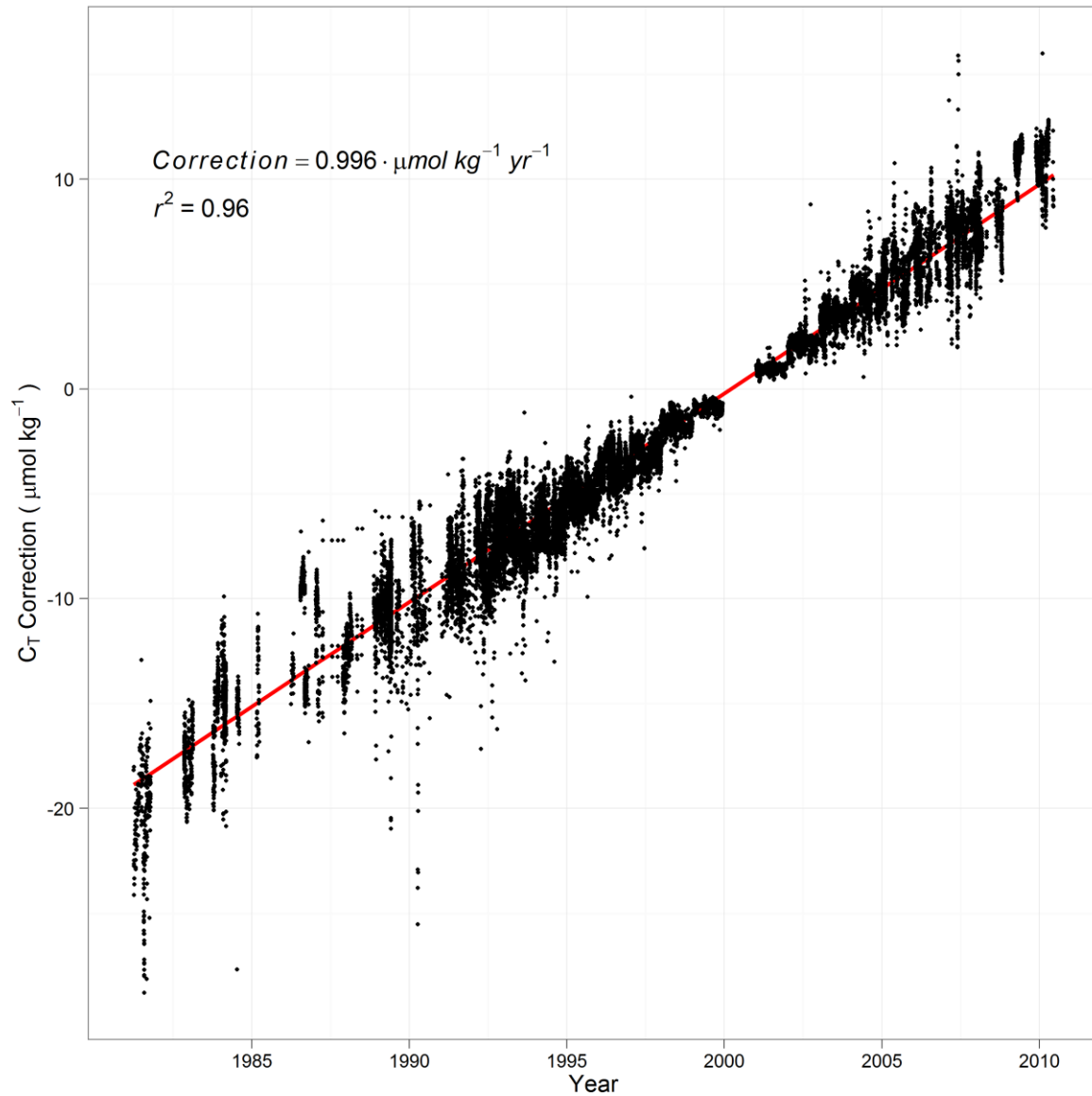
	Mean RSE ( $\mu\text{mol kg}^{-1}$ )	1 <sup>st</sup> Standard Deviation ( $\mu\text{mol kg}^{-1}$ )	% of mean
C <sub>T</sub>	12.2	0.2	1.6%
A <sub>T</sub>	8.2	0.1	1.2%

2



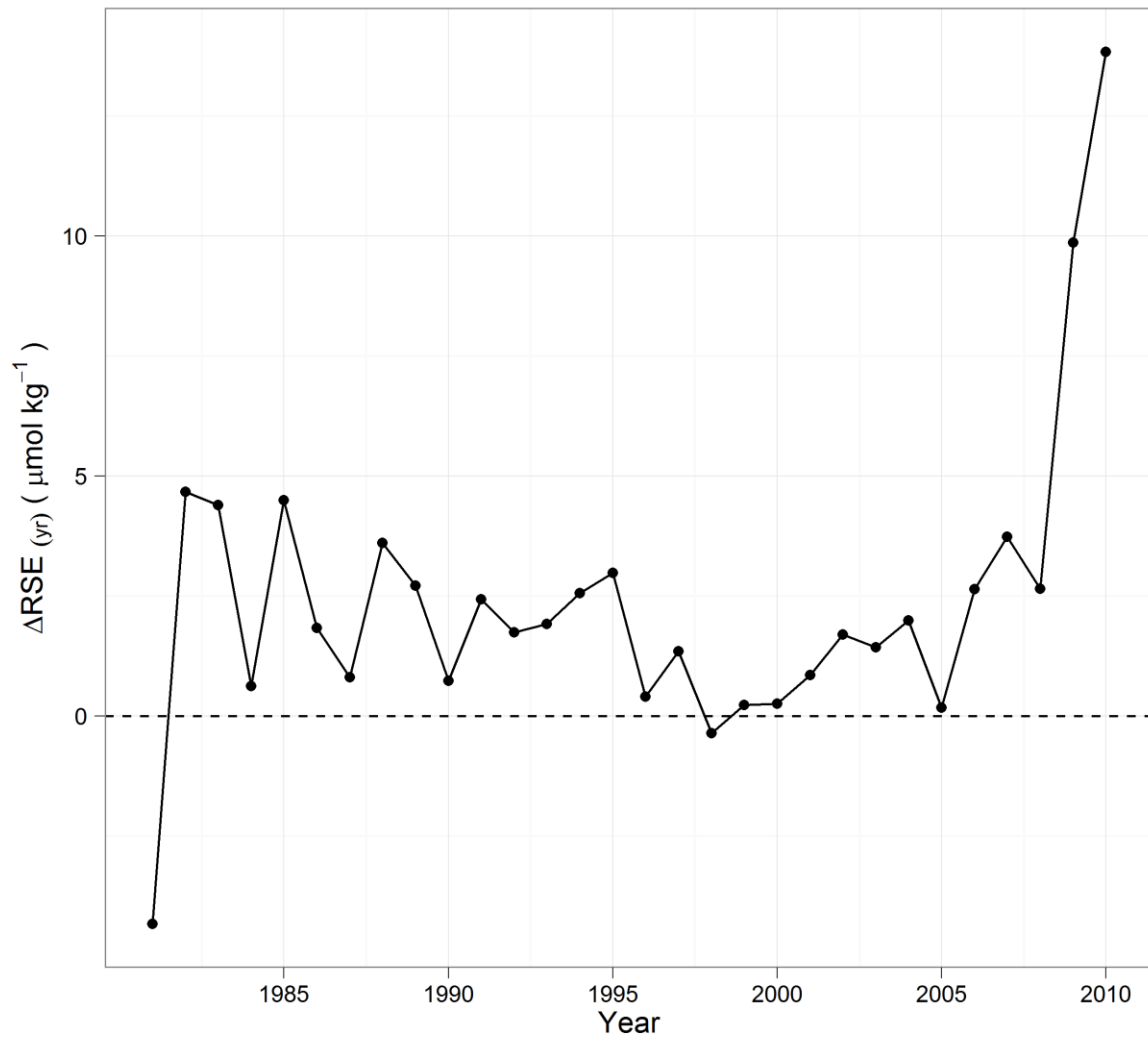
1

2 **Fig. C1.** Global and Mauna Loa site CO<sub>2</sub> difference between in-situ year and the year 2000.



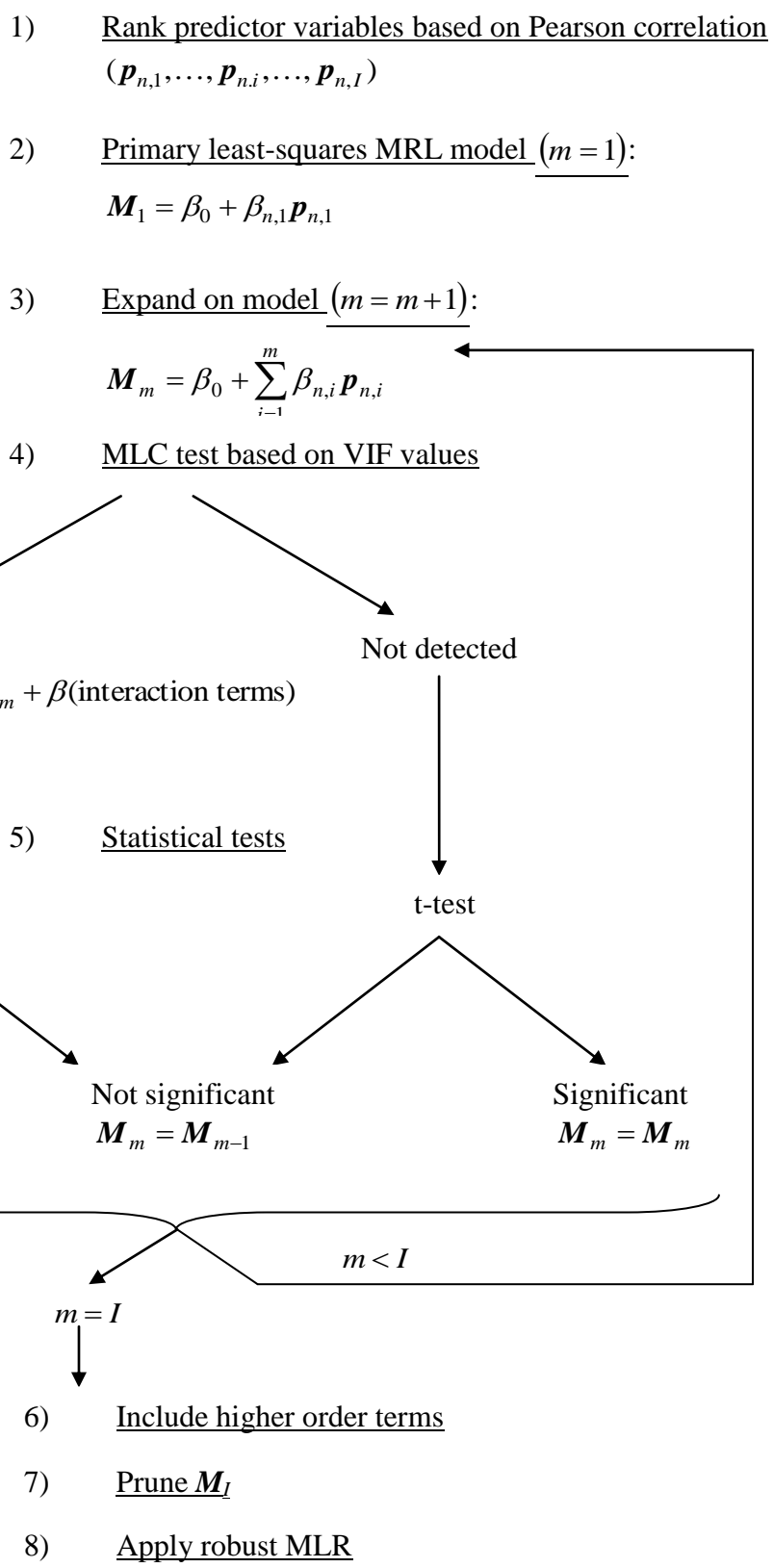
1  
2  
3  
4  
5

**Fig. C2.** Correction factor applied to  $C_T$  measurements defined by  $C_T\text{correction} = C_{T(2000)} - C_{T(\text{in-situ year})}$

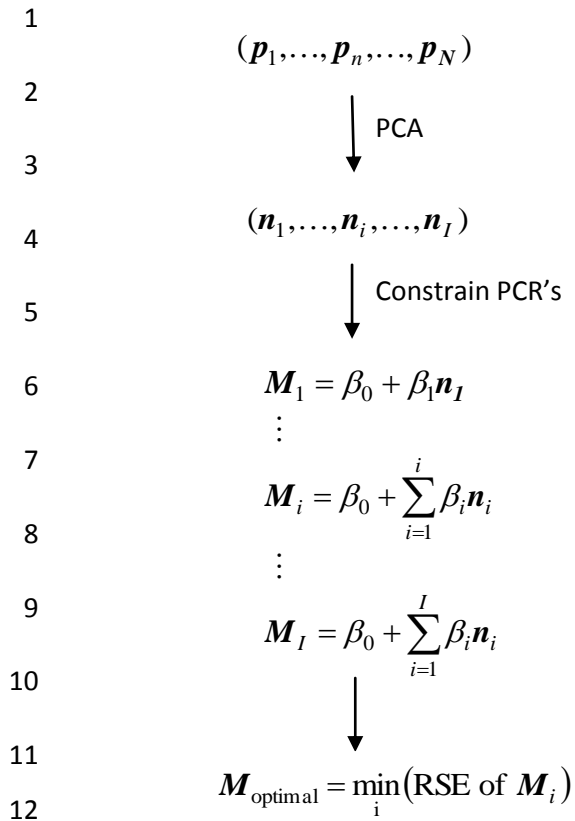


1  
2 **Fig. D1.** Annual  $\Delta RSE$  between  $C_T$  models trained with and without anthropogenic  
3 corrections.  
4

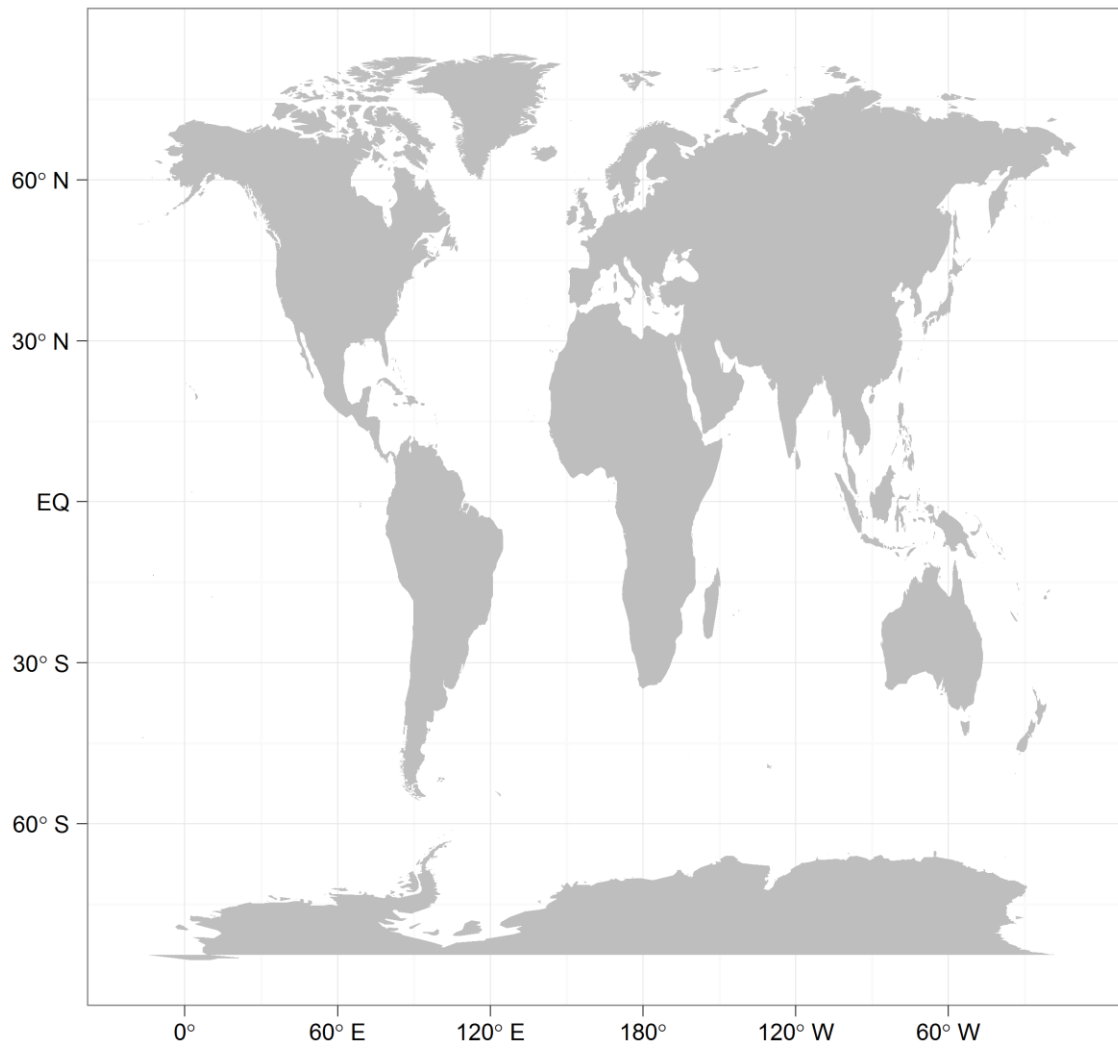
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27



**Fig. E1.** Robust forward MLR routine schematic.



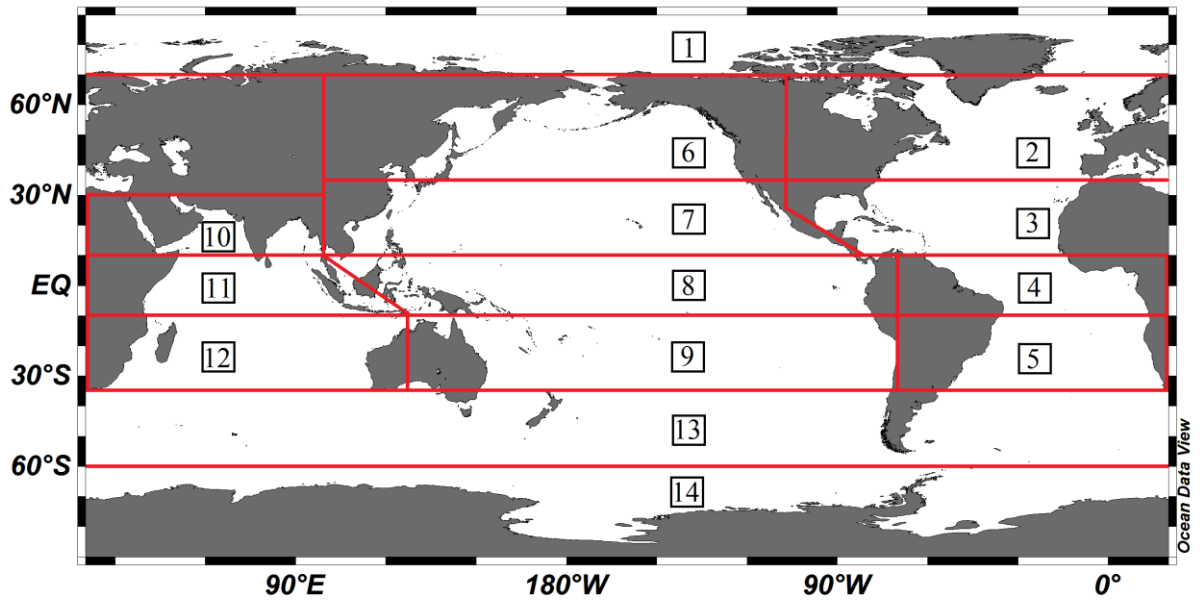
**Fig. F1.** Principle Component Regression schematic.



1  
2  
3  
4  
5

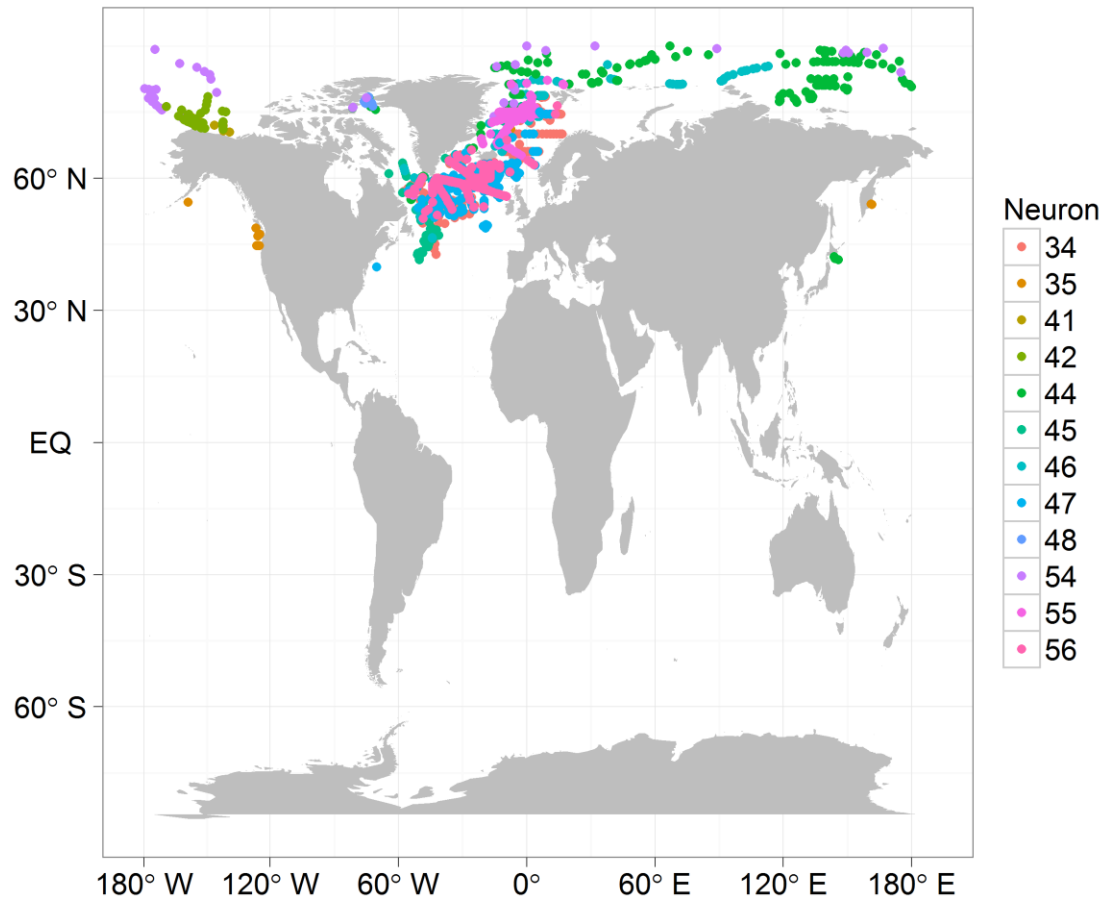
**Fig. I1.** Longitude values after reorganization.





1  
2  
3  
4

**Fig. J1.** Geographical separation of independently predicted dataset into 14 regions.



1

2 **Fig. M1.** Distribution of measurements assigned to a neuron containing at least one Arctic  
 3 measurement.

4