

1                   **Using atmospheric observations to evaluate the**  
2                   **spatiotemporal variability of CO<sub>2</sub> fluxes simulated by**  
3                   **terrestrial biospheric models**

4                   **Yuanyuan Fang<sup>1</sup>, Anna M. Michalak<sup>1</sup>, Yoichi P. Shiga<sup>1,2</sup>, Vineet Yadav<sup>1</sup>**

5                   <sup>1</sup> Department of Global Ecology, Carnegie Institution for Science, Stanford, CA, USA

6                   <sup>2</sup> Department of Civil and Environmental Engineering, Stanford University, Stanford, CA, USA

7  
8                   Correspondence to: Yuanyuan Fang (yyfang@stanford.edu)

9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29

## Abstract

Terrestrial biospheric models (TBMs) are used to extrapolate local observations and process-level understanding of land-atmosphere carbon exchange to larger regions, and serve as a predictive tool for examining carbon-climate interactions. Understanding the performance of TBMs is thus crucial to the carbon cycle and climate science. In this study, we present and assess an approach for evaluating the spatiotemporal patterns, rather than aggregated magnitudes, of net ecosystem exchange (NEE) simulated by TBMs using atmospheric CO<sub>2</sub> measurements. The approach is based on statistical model selection implemented within a high-resolution atmospheric inverse model. Using synthetic data experiments, we find that current atmospheric observations are sensitive to the underlying spatiotemporal flux variability at sub-biome scales for a large portion of North America, and that atmospheric observations can therefore be used to evaluate simulated spatiotemporal flux patterns as well as to differentiate among multiple competing TBMs. Experiments using real atmospheric observations and four prototypical TBMs further confirm the applicability of the method, and demonstrate that the performance of TBMs in simulating the spatiotemporal patterns of NEE varies substantially across seasons, with best performance during the growing season and more limited skill during transition seasons. This result is consistent with previous work showing that the ability of TBMs to model flux magnitudes is also seasonally-dependent. Overall, the proposed approach provides a new avenue for evaluating TBM performance based on sub-biome scale flux patterns, presenting an opportunity for assessing and informing model development using atmospheric observations.

## 30 **1 Introduction**

31 A key question in carbon cycle science is how terrestrial carbon sinks will evolve within the  
32 context of a rapidly changing climate. Such projections of future carbon-climate interactions  
33 largely depend on the accuracy of current terrestrial biospheric models (TBMs), the main tool  
34 used to simulate the processes controlling the biospheric carbon cycle. Thus, understanding and  
35 evaluating the performance of current TBMs is an essential step toward improving the state of  
36 carbon cycle research.

37 TBM predictions of carbon flux can be directly evaluated against eddy covariance tower  
38 measurements at various time scales ranging from hourly to interannual (Baker et al., 2003;  
39 Balzarolo et al., 2014; Keenan et al., 2012; Raczka et al., 2013; Richardson et al., 2012; Sasai et  
40 al., 2005; Schaefer et al., 2012; Schwalm et al., 2010), but the information provided by flux  
41 towers is only representative of small spatial scales ( $\sim 1\text{km}^2$ ) relative to the scales of interest for  
42 global analyses. On the other end of the spectrum, TBM predictions aggregated to large spatial  
43 and/or temporal scales (*e.g.*, continental/monthly to global/annual) are routinely intercompared  
44 with flux estimates obtained from inverse-modeling based on observed atmospheric CO<sub>2</sub> mixing  
45 ratios (Canadell et al., 2011; Gourdji et al., 2012; Hayes et al., 2012; McGuire et al., 2012;  
46 Turner et al., 2011), but such large-scale comparisons make it difficult to provide directly usable  
47 information regarding the processes driving carbon exchange. In addition, differences among  
48 TBMs exist across a full range of spatiotemporal scales, including inter-annual variability, the  
49 timing of phenology, and the spatiotemporal distribution of biospheric carbon fluxes within  
50 regions (Gourdji et al., 2012; Huntzinger et al., 2012; Keenan et al., 2012; Raczka et al., 2013;  
51 Richardson et al., 2012; Schaefer et al., 2012; Schwalm et al., 2010). These differences reflect

52 the fact that processes controlling carbon-climate feedbacks are manifested differently across  
53 TBMs.

54 Assessing the spatial and/or temporal variability of carbon fluxes as a method for evaluating  
55 TBMs, therefore, offers the potential to examine the environmental processes driving carbon  
56 exchange, and hence provides an alternative path forward in the assessment of TBM predictions.  
57 For example, evaluating the timing of modeled phenology can highlight issues associated with a  
58 model's representation of Light Use Efficiency (LUE), temperature response, and GPP response  
59 under various conditions (Richardson et al., 2012; Schwalm et al., 2010). Examining the  
60 interannual variability of TBM output can identify problems with the representation of  
61 interannual variability in spring phenology, soil thaw, snowpack melt and lagged response to  
62 extreme climatic events (Keenan et al., 2012).

63 The majority of previous studies examining carbon flux variability are still based on spatially  
64 and/or temporally aggregated carbon fluxes, however. An evaluation of flux variability, or flux  
65 patterns, at the fine native spatiotemporal scales of TBM simulations would make it possible to  
66 more directly target the fine-scale spatiotemporal patterns of carbon fluxes that have been shown  
67 to directly relate to environmental/climatic factors, such as precipitation, radiation and nighttime  
68 temperature (Beer et al., 2010; Mueller et al., 2010; Yadav et al., 2010). Such evaluations could  
69 therefore inform model improvements at the process level.

70 Observations of atmospheric CO<sub>2</sub> can potentially be used to assess such fine-scale  
71 spatiotemporal flux patterns. On one hand, atmospheric CO<sub>2</sub> observations are sensitive to fine-  
72 scale NEE spatial and temporal variability (Huntzinger et al., 2011). On the other hand,  
73 variations in atmospheric CO<sub>2</sub> measurements are routinely used in inverse modeling frameworks  
74 to infer upwind sources and sinks of CO<sub>2</sub>, and recent studies suggest that atmospheric

75 observations contain information about flux patterns at spatial and temporal resolutions  
76 comparable to those of TBMs run for regional to continental domains (Broquet et al., 2013;  
77 Göckede et al., 2010; Gourdji et al., 2010; Gourdji et al., 2012). Despite the uncertainties  
78 existing in regional inversions due to uncertainties in atmospheric transport, fossil fuel  
79 emissions, fire disturbance, and boundary conditions, these studies do point to the possibility of  
80 evaluating the spatiotemporal patterns of fluxes from biospheric models through the use of  
81 inverse models.

82 With this goal in mind, what is needed is an atmospheric-inversion-based method that can use  
83 variations in atmospheric CO<sub>2</sub> to assess the spatiotemporal patterns of surface carbon fluxes  
84 simulated by TBMs. The purpose of this paper is to present, evaluate, and demonstrate the  
85 application of such an approach, applied here to the evaluation of the 1°×1° and 3-hourly  
86 spatiotemporal variability of Net Ecosystem Exchange (NEE) simulated by TBMs using  
87 atmospheric CO<sub>2</sub> measurements. This fine scale variability is evaluated here within each month  
88 and biome over North America, thus providing a way to evaluate the seasonal and biome-  
89 specific differences in model performance. The distinguishing feature of the proposed approach  
90 is that it targets the evaluation of flux patterns at fine scales, rather than flux magnitudes at  
91 aggregated scales, thereby providing a closer link to process-based understanding of TBM  
92 performance. The approach is first evaluated with a series of synthetic data experiments where  
93 the underlying flux patterns affecting the atmospheric CO<sub>2</sub> signals are known. The application of  
94 this approach is further tested and demonstrated using actual atmospheric measurements and a  
95 prototypical small set of extensively studied TBM simulations from the North American Carbon  
96 Program (NACP) Regional Interim Synthesis (RIS) effort (Huntzinger et al., 2012).

## 97 **2 Data description**

### 98 **2.1 Atmospheric CO<sub>2</sub> measurements**

99 We use continuous, high-precision atmospheric CO<sub>2</sub> concentration measurements from 35  
100 towers for the year 2008 to evaluate the simulated NEE spatiotemporal variability over North  
101 American land. Figure 1 shows the location of these towers along with the geographic coverage  
102 of seven North American biomes as modified from Olson et al. (2001). A majority of towers are  
103 located in Temperate Broadleaf and Mixed Forests, Temperature Grasslands, Savannas and  
104 Shrublands, Temperature Coniferous Forests and Boreal Forests and Taiga, while very few  
105 towers are located in the other biomes (Tundra, Desserts and Xeric, and Tropical and Subtropical  
106 biomes). This distribution of towers is expected to affect the sensitivity of atmospheric CO<sub>2</sub> data  
107 to NEE within those biomes. The year 2008 is used as it includes the expansion of continuous  
108 measurement locations from the Mid-Continent Intensive (MCI) project (Miles et al., 2012; Ogle  
109 et al., 2006). Atmospheric CO<sub>2</sub> measurements are processed and averaged to 3-hourly intervals  
110 as described in Gourdjji et al. (2012). Data from all hours of the day are used for tall towers with  
111 a height over 300m while afternoon data are used for most short towers (lower than 100m) and  
112 nighttime data are used for sites with complex topography (e.g. Niwot Ridge - NWR), as detailed  
113 in Shiga et al. (2014). We further remove data that are strongly influenced by only a few 1°×1°  
114 grid cells, in order to exclude data that are likely subject to systematic transport model errors  
115 (Göckede et al., 2010; Gourdjji et al., 2012; Peters et al., 2007). The total number of resulting  
116 observations is  $n = 28,717$ .

117 To remove the effect of boundary conditions, we pre-subtract the GLOBALVIEW-CO2  
118 boundary condition (GLOBALVIEW-CO2, 2010) from atmospheric measurements as in Gourdjji  
119 et al. (2012). We further remove the impact of fossil fuel emissions by pre-subtracting

120 concentrations modeled based on the VULCAN-ODIAC fossil fuel emissions inventory (Shiga  
121 et al., 2014).

## 122 **2.2 Sensitivity footprints from atmospheric transport model**

123 The sensitivity of the available atmospheric observations to underlying CO<sub>2</sub> fluxes (in units of  
124 ppmv/( $\mu\text{mol m}^{-2}\text{s}^{-1}$ )) is quantified as described in Gourджи et al. (2012). In brief, footprints are  
125 derived from the Stochastic Time-Inverted Lagrangian Transport (STILT) model (Lin et al.,  
126 2003), driven by meteorological fields from the Weather Research and Forecast (WRF) model  
127 (Skamarock and Klemp, 2008). The STILT transport model has been used and examined  
128 extensively at regional and continental scales (Chatterjee et al., 2012; Gourджи et al., 2010;  
129 Gourджи et al., 2012; Huntzinger et al., 2011; Kort et al., 2008; McKain et al., 2012). Footprints  
130 are also used to generate synthetic observational time series based on TBM flux simulations.

## 131 **2.3 Terrestrial Biospheric Models (TBMs)**

132 We use simulations from four TBMs to evaluate the proposed approach, namely CASA-GFED  
133 (van der Werf et al., 2006), SiB3 (Baker et al., 2008), ORCHIDEE (Krinner et al., 2005) and  
134 VEGAS2 (Zeng et al., 2005), using the runs submitted to the NACP RIS activity. These four  
135 models were selected for analysis because of the availability of 3-hourly NEE flux output. While  
136 CASA-GFED and VEGAS2 have a coarser native temporal resolution, their NEE fluxes have  
137 been downscaled to a 3-hourly resolution as described in Huntzinger et al. (2011). Our  
138 evaluation is based on the overall NEE simulated by each TBM, although model definitions of  
139 NEE differ: CASA-GFED includes fire disturbance while other models do not; ORCHIDEE  
140 exclude crop harvest while others do not. A comparison and summary of these simulations can  
141 be found in Table S1 in the supplementary material. Further details on the NACP RIS  
142 simulations can be found in Huntzinger et al. (2012).

### 143 **3 Regression framework linking atmospheric CO<sub>2</sub> to NEE**

144 The overall goal of the proposed approach is to evaluate the spatiotemporal variability of NEE as  
145 simulated by various TBMs using atmospheric CO<sub>2</sub> measurements. Such an approach must be  
146 based on an inverse model that can infer NEE from atmospheric CO<sub>2</sub> measurements. It must  
147 also include a statistical model selection component to evaluate the degree to which NEE  
148 patterns predicted by TBMs are useful in explaining the observed atmospheric CO<sub>2</sub> variability.  
149 Rather than quantifying the magnitude of NEE, the primary goal here is to evaluate the  
150 spatiotemporal NEE patterns (at a 1°×1° and 3-hourly resolution) within specific biomes of  
151 North America and for specific months. The approach presented here builds on the geostatistical  
152 inverse modeling (GIM) framework (Gourdji et al., 2010; Gourdji et al., 2012; Michalak et al.,  
153 2004), but is presented here in the form of a regression analysis to simplify the presentation and  
154 emphasize the introduction of model selection aspect of the proposed approach.

155 To this end, we first formulate a multi-linear regression framework that relates atmospheric  
156 observations to NEE spatiotemporal variability. Statistical model selection is then applied to  
157 determine whether, when, and where the spatiotemporal variability of simulated NEE is  
158 consistent with that evident from variability in atmospheric CO<sub>2</sub>. Here, the NEE spatiotemporal  
159 variability is defined at a 1°×1° spatial and 3-hourly temporal resolution, and the TBMs are  
160 evaluated within specific biome-month combinations. Figure 2 shows the distribution of NEE in  
161 one specific biome-month combination (*i.e.*, Boreal Forests and Taiga in July) as an example.

162 To link atmospheric measurement to surface fluxes we first define the observed atmospheric CO<sub>2</sub>  
163 concentrations, with the influence of boundary conditions and fossil fuel emissions pre-  
164 subtracted, as:



$$\mathbf{z} = \mathbf{H}\mathbf{s} + \boldsymbol{\varepsilon} \quad (1)$$

165 where  $\mathbf{z}$  is an  $n \times 1$  vector of atmospheric CO<sub>2</sub> observations,  $\mathbf{s}$  is an  $m \times 1$  vector of the  
 166 underlying NEE fluxes at  $1^\circ \times 1^\circ$  and 3-hourly resolution,  $\mathbf{H}$  ( $n \times m$ ) are the sensitivity  
 167 footprints, namely a Jacobian matrix representing the sensitivity of each observation to each  
 168 underlying flux (i.e.,  $\frac{\partial z_i}{\partial s_j}$ ) as quantified using an atmospheric transport model (see Section 2.2),  
 169 and  $\boldsymbol{\varepsilon}$  ( $n \times 1$ ) is the model-data mismatch term that represents any discrepancies between  
 170 observed ( $\mathbf{z}$ ) and modeled ( $\mathbf{H}\mathbf{s}$ ) CO<sub>2</sub> mixing ratios. The model-data mismatch term encompasses  
 171 the influence of errors in the boundary conditions, errors in the fossil fuel inventory,  
 172 representation errors, aggregation errors, transport model errors, and measurement errors. These  
 173 errors are assumed to have zero mean and be uncorrelated across measurements, with their  
 174 variances represented by a diagonal covariance matrix  $\mathbf{R}$  ( $n \times n$ ). The dimensions of the  
 175 matrices and vectors are based on the total number of observations,  $n = 28,717$ , and the total  
 176 number of fluxes at a  $1^\circ \times 1^\circ$  (2635 such grid cells within the domain used here) and 3-hourly  
 177 resolution ( $366 \times 8 = 2928$ ) such periods within the span of the one-year inversion),  $m =$   
 178  $2635 \times 2928 = 7,715,280$ .

179 The spatiotemporal NEE distribution of  $\mathbf{s}$  is represented as a linear model of NEE as predicted  
 180 by various TBMs within specific biome-month combinations:

$$\mathbf{s} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\xi} \quad (2)$$

181 where  $\mathbf{X}$  is a  $m \times p$  matrix with each column representing NEE  $1^\circ \times 1^\circ$  3-hourly spatiotemporal  
 182 variability within a specific biome-month combination from a specific TBM, such that a given  
 183 column is populated by the modeled NEE from a given TBM for a given biome-month for those  
 184 rows (i.e. elements of  $\mathbf{s}$ ) corresponding to that specific biome-month combination, while the

185 remainder of the column is filled with zeros. These individual columns of  $\mathbf{X}$  are thus predictor  
 186 variables for the dependent variable  $\mathbf{s}$ . With 7 biomes (Figure 1) and 12 months, there are a total  
 187 of 84 possible predictor variables for each TBM (i.e.,  $p \leq 84$  for one TBM). The  $p \times 1$  vector  $\boldsymbol{\beta}$   
 188 represents the drift coefficient describing the relationship between  $\mathbf{X}$  and  $\mathbf{s}$ , and  $\mathbf{X}\boldsymbol{\beta}$  together thus  
 189 represents a statistical model of the trend of NEE. The  $m \times 1$  vector  $\boldsymbol{\xi}$  represents the portion of  
 190 the variability of  $\mathbf{s}$  that cannot be explained by the predictor variables in  $\mathbf{X}$ , and these deviations  
 191 are modeled as having a mean of zero and a covariance matrix  $\mathbf{Q}$  ( $m \times m$ ) that represents how  
 192 the flux deviations from the model of the trend (i.e.,  $\mathbf{s} - \mathbf{X}\boldsymbol{\beta}$ ) are correlated in time and space.  
 193 Combining these two equations, we represent the atmospheric observations  $\mathbf{z}$  in terms of the  
 194 NEE predictor variables  $\mathbf{X}$ :

$$\mathbf{z} = \mathbf{H}\mathbf{X}\boldsymbol{\beta} + \mathbf{H}\boldsymbol{\xi} + \boldsymbol{\varepsilon} \quad (3)$$

195 where  $\mathbf{z}$  is seen to have a spatiotemporally variable mean  $\mathbf{H}\mathbf{X}\boldsymbol{\beta}$  and, assuming independence  
 196 between  $\boldsymbol{\xi}$  and  $\boldsymbol{\varepsilon}$ , a residual covariance of:

$$\boldsymbol{\Sigma} = \mathbf{H}\mathbf{Q}\mathbf{H}^T + \mathbf{R} \quad (4)$$

197 where  $T$  is the matrix transpose operation. From a statistical standpoint, our goal then becomes  
 198 to select a subset of TBM biome-month combinations that capture a substantial portion of the  
 199  $\text{CO}_2$  variability observed in  $\mathbf{z}$ . This constitutes a classical statistical model selection problem, in  
 200 which we examine which predictor variables (candidate columns in  $\mathbf{X}$ ) are useful in explaining  
 201 the atmospheric  $\text{CO}_2$  measurements ( $\mathbf{z}$ ).

202 A widely applied approach for statistical model selection is the Bayesian Information Criterion  
 203 (BIC) (Schwarz, 1978). BIC takes into account both the goodness of fit, i.e., the residual sum of  
 204 squares ( $RSS$ ), and the numbers of auxiliary variables ( $k$ ) in each candidate model, and can be

205 used to compare non-nested candidate models. BIC has also been adapted for use with  
 206 spatiotemporally autocorrelated residuals (Hoeting et al., 2006; Mueller et al., 2010) and within  
 207 the context of atmospheric inversions where atmospheric observations are used to inform  
 208 underlying surface fluxes (Gourdji et al., 2012), making it ideal for the application presented  
 209 here. The standard expression for BIC is:

$$BIC = \underbrace{\ln|\boldsymbol{\Sigma}| + RSS}_{\text{log likelihood}} + \underbrace{k \ln(n)}_{\text{penalty term}} \quad (5)$$

210 where  $RSS$  represents the residual sums of squares of a given candidate model  $\mathbf{X}_c$ ,  $\boldsymbol{\Sigma}$  is the  $n \times n$   
 211 covariance matrix of the residuals (Eq. 4),  $||$  denotes the matrix determinant, and  $k$  is the number  
 212 of parameters in a particular candidate model. For the specific application presented here (Eq. 1-  
 213 4) and factoring out the unknown drift coefficients,  $\boldsymbol{\beta}$  and  $RSS$  become as in Gourdji et al.  
 214 (2012):

$$\boldsymbol{\beta} = ((\mathbf{H}\mathbf{X}_c)^T \boldsymbol{\Sigma}^{-1} (\mathbf{H}\mathbf{X}_c))^{-1} (\mathbf{H}\mathbf{X}_c)^T \boldsymbol{\Sigma}^{-1} \mathbf{z} \quad (6)$$

$$RSS = \left[ \mathbf{z}^T \left( \boldsymbol{\Sigma}^{-1} - \boldsymbol{\Sigma}^{-1} (\mathbf{H}\mathbf{X}_c) ((\mathbf{H}\mathbf{X}_c)^T \boldsymbol{\Sigma}^{-1} (\mathbf{H}\mathbf{X}_c))^{-1} (\mathbf{H}\mathbf{X}_c)^T \boldsymbol{\Sigma}^{-1} \right) \mathbf{z} \right] \quad (7)$$

215 The specific covariance parameters needed to define  $\mathbf{Q}$  and  $\mathbf{R}$ , which are needed to define  $\boldsymbol{\Sigma}$ , vary  
 216 between experiments and are obtained as described in the supplementary materials.

217 Model selection built on this framework aims to identify the “best” model of the trend based on a  
 218 tradeoff between model size and the model’s power in explaining the variations in observed  
 219 atmospheric CO<sub>2</sub>. Here, the “best” model is specially defined as one with the minimum BIC  
 220 value, providing an optimal balance between model complexity and model fit. To identify this  
 221 model, BIC is compared across all possible combinations of predictor variables (i.e. 84 NEE  
 222 biome-months per TBM). Due to the large number of candidate predictor variables considered

223 here, we implement the branch-and-bound algorithm of Yadav et al. (2013) to improve  
224 computational efficiency.

225 The final selected subset of TBM biome-months represents those biomes and months within  
226 which a given TBM exhibits spatiotemporal variability that explains a substantial portion of the  
227 variability observed in the observations  $\mathbf{z}$  (see Eq. 3). For a given TBM biome-month  
228 distribution to be “selected” as part of the “best” model of the trend, therefore, (1) the available  
229 atmospheric observations must be sensitive to the spatiotemporal variability of fluxes within that  
230 biome-month (as represented through  $\mathbf{H}$ ), i.e., the information contained in atmospheric data  
231 sufficiently constrains the spatiotemporal variability within that biome-month, and (2) the  
232 variability within a particular biome-month as represented by a particular TBM must explain a  
233 sufficient portion of the variability in the atmospheric observations to offset the penalty term in  
234 Eq. (5), i.e. the reduction in  $RSS$  must outweigh the penalty term. On the contrary, if a given  
235 TBM biome-month distribution is “not selected” then, either (1) or (2) as given above is not  
236 satisfied, i.e., either that atmospheric observations are not sensitive to the NEE variability within  
237 that biome-month, or that the NEE variability as represented in the model is inconsistent with  
238 atmospheric observations. In other words, selecting or not selecting a TBM biome-month  
239 combination directly reflects on the performance of the TBM in that biome and month, as long as  
240 we have fulfilled the requirement in (1) above. If the condition in (1) is not met, we are not able  
241 to use the model selection results to examine model performance, due to the insufficient  
242 coverage of the network. We henceforth refer to the TBM biome-month combinations included  
243 in the final selected subset as the “selected” combinations or elements, or alternately as the TBM  
244 biome-month combinations “identified” using the atmospheric data.

#### 245 **4 Synthetic data and real data experiments**

246 In this Section, we design a series of Synthetic Data (SD) experiments (Figure 3), in which the  
247 underlying fluxes are prescribed, to assess the sensitivity of atmospheric CO<sub>2</sub> measurements to  
248 NEE flux spatiotemporal patterns within all biome-month combinations, and identify when and  
249 where results from the proposed approach reliably reflect model performance in simulating NEE  
250 spatiotemporal variability. We further introduce two Real Data (RD) experiments as a proof-of-  
251 concept demonstration of our approach. In those RD experiments, we use actual atmospheric  
252 CO<sub>2</sub> measurements to evaluate the spatiotemporal variability of NEE as simulated by four  
253 prototypical TBMs (Section 2.3).

254 In the SD experiments, synthetic atmospheric observations ( $\mathbf{z}$ ) are generated as described in Eq.  
255 (1) using fluxes ( $\mathbf{s}_{\text{TBM}}$ ) that include NEE as simulated by one of the TBMs and, in some cases,  
256 spatiotemporally-correlated flux residuals ( $\boldsymbol{\xi}$ ) and model-data mismatch errors ( $\boldsymbol{\epsilon}$ ), *i.e.*,  $\mathbf{z} =$   
257  $\mathbf{H}(\mathbf{s}_{\text{TBM}} + \boldsymbol{\xi}) + \boldsymbol{\epsilon}$ . The superset of candidate ancillary variables (Figure 3,  $\mathbf{X}$ ) includes NEE from  
258 one or more TBMs. TBMs included in  $\mathbf{s}_{\text{TBM}}$  and  $\mathbf{X}$  are denoted as the “truth” and the “candidate  
259 (s)” respectively henceforth.

260 The first SD case study, SD-one- $\emptyset\emptyset$  (Figure 3), is designed to investigate whether, when, and  
261 where the information contained in current atmospheric data enables the identification of the  
262 correct candidate TBM for a case where it is the only TBM considered in the model selection,  
263 where this TBM fully represents the variability in the synthetic atmospheric observations ( $\boldsymbol{\xi}=0$ ),  
264 and where no model-data mismatch errors are included in the simulation ( $\boldsymbol{\epsilon}=0$ ). Given that in this  
265 case the candidate TBM explains all of the variability in the synthetic atmospheric observations,  
266 it should always be selected if the atmospheric data are sufficiently sensitive to NEE across all

267 biome-months; hence, biome-months for which the TBM is not selected are ones to which the  
268 atmospheric CO<sub>2</sub> observations are not sufficiently sensitive to offset the penalty term in Eq. (5).

269 The second and third SD case studies, SD-one- $\emptyset\epsilon$  and SD-one- $\xi\epsilon$  (Figure 3), are analogous to  
270 SD-one- $\emptyset\emptyset$ , but include model-data mismatch errors ( $\epsilon \neq 0$ , denoted by  $\epsilon$ ) and/or spatially  
271 correlated flux residuals ( $\xi \neq 0$ , denoted by  $\xi$ ). These case studies are designed to test the degree  
272 to which current atmospheric observations can inform the spatiotemporal variability of NEE in  
273 cases with realistic model-data mismatch errors and where the candidate TBM only represents a  
274 portion of the true underlying NEE variability. In these case studies, noise ( $\epsilon$ ) is added to  
275 observations, generated as a random vector of independent normally-distributed values with  
276 variances corresponding to the diagonal elements of  $\mathbf{R}$ , which are inferred from the RD-all- $\xi\epsilon$   
277 experiment (described below), and a mean of 0. In addition for SD-one- $\xi\epsilon$ , the flux signal from  
278 the TBMs is augmented with additional spatially-correlated fluxes ( $\xi$ ) generated as a random  
279 vector of normally distributed values with a covariance structure equal to that inferred from the  
280 RD-all- $\xi\epsilon$  experiment (described below). The details of the model-data mismatch errors and flux  
281 residuals are summarized in the supplementary materials.

282 The final SD case study, SD-all- $\xi\epsilon$ , builds on SD-one- $\xi\epsilon$  (Figure 3), but is designed to test  
283 whether the correct TBM can be identified when all four TBMs are used as candidate variables.  
284 This case study therefore explores whether current atmospheric observations can be used to  
285 differentiate among candidate TBMs. No constraints are placed on the model selection, such that  
286 more than one TBM can be selected for the same biome-month, but only the dominant TBM (i.e.  
287 the one with the largest  $\beta$ , Eq. 6) is discussed in analyzing this case.

288 Finally, two RD case studies, RD-one- $\xi\epsilon$  and RD-all- $\xi\epsilon$ , are defined analogously to SD-one- $\xi\epsilon$   
289 and SD-all- $\xi\epsilon$ , to further test the applicability of our approach by examining the actual  
290 performance of the four prototypical TBMs based on available atmospheric observations. The  
291 observations ( $\mathbf{z}$ ) here are the actual atmospheric measurements, which by definition encompass  
292 model-data mismatch errors, and the flux residuals are also inherently present as no TBM is  
293 expected to perfectly reflect the true underlying fluxes. In each RD-one- $\xi\epsilon$  experiment, one of  
294 the four prototypical TBMs is used as the candidate TBM in order to assess individual TBM  
295 performance. In RD-all- $\xi\epsilon$ , all four TBMs are included, analogously to SD-all- $\xi\epsilon$ , to identify the  
296 TBM (if any) that best represents the spatiotemporal variability of NEE within a given biome-  
297 month, based on the information provided by the atmospheric measurements.

## 298 **5 Sensitivity of atmospheric observations to NEE spatiotemporal** 299 **variability and evaluation of the proposed approach**

300 The SD-one- $\emptyset\emptyset$  experiment examines the sensitivity of atmospheric observations to underlying  
301 flux variability and evaluates the proposed approach under idealized conditions where the true  
302 flux field is perfectly represented by the candidate TBM model, and where no model-data  
303 mismatch errors are included in the synthetic atmospheric observations.

304 Results indicate that the candidate TBM is selected for over 90% of all biome-months (Figure 4,  
305 top row), demonstrating that atmospheric observations are sensitive to NEE spatiotemporal  
306 variability, and that the proposed approach leverages this sensitivity to correctly identify the  
307 TBM model as being representative of the flux variability within the vast majority of biomes and  
308 months. The only notable exception is for the Tundra biome for which, other than during the  
309 height of the growing season, the atmospheric data do not provide a sufficient constraint on the  
310 flux variability, due to the poor data coverage and the weak biospheric signal. Because this  
311 biome is expected to play an important role in future global carbon cycle and climate (Belshe et

312 al., 2013; Ping et al., 2008; Schuur et al., 2009; Tarnocai et al., 2009) and large uncertainties  
313 remain in quantifying its role in carbon cycling (McGuire et al., 2012), this result highlights the  
314 need for strategic placement of additional CO<sub>2</sub> monitoring stations in the vicinity of this biome to  
315 constrain its carbon flux distribution.

316 The SD-one- $\emptyset\epsilon$  and SD-one- $\xi\epsilon$  case studies examine the degree to which the presence of model-  
317 data mismatch errors and additional flux variability not represented by the candidate TBM limit  
318 the information content of available observations, and the ability of the proposed approach to  
319 identify the consistency between the true underlying NEE patterns and those simulated by  
320 TBMs.

321 Results of SD-one- $\emptyset\epsilon$  show that including realistic model-data mismatch errors decreases the  
322 information content of atmospheric observations to the point where a TBM that in reality  
323 represents the full spatiotemporal flux variability is not selected for many months and TBMs  
324 within the Tropical and Subtropical biome, as well as the Desert and Xeric Shrublands biome, in  
325 addition to the Tundra biome that was not well constrained even under idealized conditions  
326 (Figure 4, middle row). The identification of a TBM as correctly representing the flux patterns  
327 also becomes more challenging during winter and spring within the Boreal Forests and Taiga  
328 biome, and the Temperate Coniferous Forests biome (Figure 4, middle row), especially when  
329 VEGAS2 is used as the true flux distribution. This result is related to the fact that the magnitude  
330 and the spatiotemporal variability of NEE simulated by VEGAS2 within those biome-months are  
331 much smaller than for other TBMs. For example, the standard deviation of NEE simulated by  
332 VEGAS2 is less than a half of that of other TBMs. Overall, the inclusion of realistic model-data  
333 mismatch, combined with the coverage of the monitoring network, make the identification of  
334 TBMs that represent the spatiotemporal variability of fluxes within biomes unreliable for three of



335 the seven biomes considered here, namely the Tundra, Tropical and Subtropical, and Desert and  
336 Xeric Shrublands biomes. Subsequent analyses therefore focus on the remaining four better-  
337 constrained biomes, namely the (i) Boreal Forests and Taiga, (ii) Temperate Coniferous Forests,  
338 (iii) Temperate Grasslands, Savannas, and Shrublands, and (iv) Temperate Broadleaf and Mixed  
339 Forests biomes.

340 SD-one- $\xi\epsilon$  is the most realistic single-TBM synthetic data experiment, as it includes not only  
341 model-data mismatch errors, but also variability in the spatiotemporal flux distribution that is not  
342 represented by the candidate TBM. Results for the better-constrained biomes indicate that the  
343 ability to identify a model as correctly representing a portion of the true flux variability  
344 deteriorates in the winter months for the Boreal Forests and Taiga, but remains largely  
345 unchanged in the other biomes (Figure 4, bottom row). For the winter in the Boreal Forests and  
346 Taiga biome, the TBM is only identified when the fluxes are based on SiB3, likely because this  
347 TBM has a stronger flux signal in this biome during the winter relative to the other TBMs,  
348 thereby overcoming the confounding impacts of model-data mismatch errors and additional flux  
349 variability unexplained by the TBM.

350 Results of SD-one- $\xi\epsilon$  indicate that under realistic conditions, the proposed approach is able to  
351 correctly identify a TBM that represents a portion of the true underlying flux variability within  
352 four of the seven biomes considered here, given the monitoring network used here. The  
353 magnitude of the model data mismatch used here was derived from the real-data experiments  
354 (RD-one- $\xi\epsilon$ ), and includes the impact of errors in the transport model, boundary conditions,  
355 fossil fuel emissions, and fire emissions, as well as measurement and aggregation errors.  
356 Therefore, results suggest that conclusions over the four considered biomes are robust in spite of  
357 the influences of those uncertainties. We acknowledge that the errors applied do not fully address

358 the complexity of uncertainties in the real world, as we assume that errors are independent and  
359 follow a Gaussian distribution. However, the results presented here, together with evidence from  
360 the literature (*e.g.*, Gourdji et al., 2012; Pillai et al., 2012), support the ability to infer flux  
361 patterns despite the many sources of uncertainty in regional inversions.

362 The final SD case, SD-all- $\xi\epsilon$ , is designed to explore whether atmospheric observations can be  
363 used to differentiate among several competing TBMs to identify the TBM that best represents the  
364 underlying flux variability. Results indicate that across the majority of the examined biomes,  
365 months, and TBMs, the proposed approach combined with the available atmospheric data are  
366 able to discriminate among models for a similar fraction of TBM-biome-month combination  
367 (Figure 5) as when only the “correct” TBM was offered as a candidate model (SD-one- $\xi\epsilon$ , Figure  
368 4, bottom row). One noticeable difference, however, occurs during the growing season in the  
369 Boreal Forests and Taiga when VEGAS2 or CASA-GFED is used to represent a substantial  
370 portion of the true flux variability. In these cases, the other of these two models is often  
371 identified in the model selection procedure. This is not surprising, because these two models  
372 yield fluxes that are highly spatiotemporally correlated to one another (Figure 6), and because  
373 biospheric signals simulated by VEGAS2 are particularly weak (Huntzinger et al., 2011).  
374 Overall, therefore, for the four better-constrained biomes, the information content of the  
375 atmospheric data is sufficient to identify a TBM that represents a substantial portion of the true  
376 underlying variability using the proposed approach, even when multiple competing TBMs are  
377 available. In other words, atmospheric observations can be used to differentiate among  
378 competing TBMs. The exception, not surprisingly, is when the competing TBMs have fluxes  
379 that are highly correlated ( $R > 0.8$ ), which, for the four TBMs examined here, occurs most often  
380 over the Boreal Forests and Taiga and Temperate Coniferous Forests biomes (where biospheric

381 signals are relative weak and atmospheric data are less sensitive), for the VEGAS2 and CASA-  
382 GFED as well as SiB3 and ORCHIDEE model pairs (Figure 6).

## 383 **6 Demonstration of the proposed approach using atmospheric** 384 **observations**

385 The results presented in Section 5 confirm that, given the coverage of atmospheric data available  
386 in 2008, the proposed approach is able to identify TBMs representing a substantial portion of the  
387 underlying NEE spatiotemporal variability over four better-constrained biomes of North America  
388 throughout most of the year. In this Section, by focusing on the RD experiment results, we  
389 demonstrate the application of the proposed approach using “real” data, by evaluating four  
390 prototypical TBMs participating in the NACP RIS.

### 391 **6.1 Performance of TBMs in simulating the spatiotemporal variability of NEE**

392 The RD-one- $\xi\epsilon$  case study includes four experiments, each evaluating one prototypical TBM.  
393 As a general indication of individual TBM performance across biomes and months, we sum the  
394 number of candidate TBMs selected across the four RD-one- $\xi\epsilon$  cases (Figure 7). We find that  
395 the capability of TBMs to simulate the NEE spatiotemporal variability varies strongly across  
396 biomes and seasons. TBMs are most frequently identified over the Temperate Broadleaf and  
397 Mixed Forests biome (7 out of 12 months with at least one TBM identified), and least frequently  
398 identified over the Boreal Forests and Taiga biome (3 out of 12 months with at least one TBM  
399 identified). Across seasons, TBMs are most frequently identified during the growing season  
400 (May-Sept, 15 out of 20 biome-months with at least one TBM identified). TBMs are least  
401 frequently identified during transition seasons (Mar-Apr and Oct-Nov, with 2 out of 16 biome-  
402 months with at least one TBM identified), likely reflecting known challenges of TBMs in  
403 representing the seasonal cycle of phenology (Richardson et al., 2012; Schaefer et al., 2012;  
404 Schwalm et al., 2010). Specifically, during Oct-Nov, none of the TBMs is identified as

405 representing the flux spatiotemporal variability in any of the biomes, in agreement with the  
406 finding in Gourdji et al. (2012) that carbon fluxes simulated by over 70% of the NACP TBMs  
407 are outside the 95% confidence intervals of atmospheric inversion estimates in October.

408 Of all 48 biome-months examined, none of the four TBMs are identified as substantially  
409 representing the spatiotemporal variability within 27 biome-months, and only one TBM is  
410 identified in 5 additional biome-months (Figure 7). Multiple TBMs are identified as representing  
411 a portion of the spatiotemporal variability within the remaining 16 biome-months (Figure 7).  
412 Interestingly, SiB3 and ORCHIDEE are selected in almost all of these 16 biome-months,  
413 suggesting that they both have the potential to explain a substantial portion of the observed  
414 variability in atmospheric CO<sub>2</sub>. This is consistent with the similarity in NEE spatiotemporal  
415 series between SiB3 and ORCHIDEE shown in Figure 6.

416 The RD-all- $\xi\epsilon$  case study identifies the TBM that best represents the underlying flux variability  
417 (Figure 8). Out of 27 biome-months for which no individual TBM was selected in the RD-one- $\xi\epsilon$   
418 experiments, 5 biome-months lead to models being selected when more than one model can be  
419 used in combination, with the dominant TBM being ORCHIDEE over the Temperate Coniferous  
420 Forests biome in Apr and May and the Temperate Broadleaf and Mixed Forests in Feb, SiB3  
421 over the Boreal Forests and Taiga in Aug, and VEGAS2 over the Temperate Grasslands,  
422 Savannas and Shrublands in Dec.

423 Overall, SiB3 and ORCHIDEE are selected as the dominant TBM in explaining the flux  
424 variability as observed through the atmospheric CO<sub>2</sub> measurements more often than VEGAS2  
425 and CASA-GFED (Figure 8). SiB3 appears most representative of flux patterns over boreal  
426 biomes, whereas ORCHIDEE is most representative over temperate biomes. Although SiB3

427 appears to be selected most often (13 biome-months), followed by ORCHIDEE (10 biome-  
428 months), none of the TBMs is consistently better than the others across all biomes and seasons.

## 429 **6.2 Evaluation of the TBMs and the proposed approach within the context of earlier** 430 **studies**

431 To further evaluate the performance of, and value added provided by, the proposed approach, we  
432 assess the RD-one- $\xi\epsilon$  results within the context of the existing literature to determine whether  
433 (1) results are consistent with the literature wherever they are comparable, and (2) the proposed  
434 approach can provide insights that go beyond those provided by other model evaluation  
435 strategies.

436 Many of our findings are consistent with early work analyzing the examined TBMs within the  
437 framework of the NACP RIS. For example, we find distinctive seasonal differences in TBM  
438 performance in simulating NEE (Figures 7 and 8), consistent with the previously noted model  
439 misrepresentation of phenology seasonality based on site-level measurements (Richardson et al.,  
440 2012; Schaefer et al., 2012; Schwalm et al., 2010). In addition, we find that models perform  
441 better for Temperate Broadleaf and Mixed Forests, and that SiB3 appears to be more consistent  
442 with observations than other models, both of which are consistent with existing literature  
443 evaluating NACP RIS models (Raczka et al., 2013; Schwalm et al., 2010). The consistency  
444 between our results and existing literature further supports the performance of the proposed  
445 approach. It also implies that, although the approach proposed here is subject to many of the  
446 same uncertainties in fossil fuel emissions, fire disturbance, boundary conditions and transport  
447 models that affect all regional inversions, the main conclusions regarding TBM performance for  
448 the four major biomes examined here are quite robust.

449 The proposed approach also provides the opportunity to draw conclusions that go beyond the  
450 current literature. We present two examples here.

451 First, results indicate that model capability in simulating the spatiotemporal variability (i.e.  
452 patterns) of NEE varies strongly with seasons, with greater skill during the growing season than  
453 during the transition seasons. In other words, even within specific biomes and months, the  
454 variability of NEE is better represented during the growing season. This seasonal variability in  
455 model performance may be due to seasonal differences in the dominant environmental drivers  
456 controlling the spatiotemporal variability of NEE. For example, Mueller et al. (2010) found that  
457 the environmental drivers controlling NEE at a hardwood forest vary across seasons, with  
458 radiation, nighttime temperature and vegetative radiation indices (*i.e.*, fPAR) dominating during  
459 the growing, non-growing and leaf-out seasons, respectively. We hypothesize that the seasonal  
460 differences in model performance is likely related to the models' ability to represent the  
461 seasonally-varying influence of such environmental drivers. Because the NEE spatiotemporal  
462 variability is directly related to environmental processes and drivers (Beer et al., 2010; Mueller et  
463 al., 2010; Yadav et al., 2012; Gourdji et al., 2012), the proposed approach provides a close link  
464 between model performance and environmental processes.

465 Second, we find that SiB3 and ORCHIDEE are identified more often as representing the  
466 spatiotemporal flux variability than VEGAS2 and CASA-GFED. Given that the simulated NEE  
467 spatiotemporal variability is more similar between SiB3 and ORCHIDEE, and between VEGAS2  
468 and CASA-GFED, relative to across these two model pairs (Figure 6), this finding suggests that  
469 aspects of the model internal structure common within the pairs likely contribute to similarities  
470 in simulated flux patterns and associated performance. Such features include: 1) SiB3 and  
471 ORCHIDEE use Enzyme Kinetic (EK) models while CASA-GFED2 and VEGAS use Light Use

472 Efficiency (LUE) models to formulate their photosynthesis processes; 2) the native model time  
473 step of SiB3 and ORCHIDEE is shorter than a day while that of CASA-GFED and VEGAS2  
474 varies from daily to monthly; and 3) SiB3 and ORCHIDEE have substantially more plant  
475 functional types (PFTs) than CASA-GFED and VEGAS2. Although it is not possible to draw  
476 definite conclusions about the links between model structure and model performance in  
477 simulating flux patterns based on the small number of TBMs examined here and the lack of a  
478 uniform simulation protocol, a future application of this approach to a larger ensemble of models  
479 following a uniform protocol would make it possible to explore these connection in more detail.

## 480 **7 Concluding remarks**

481 In this paper, we present, evaluate and demonstrate a statistical approach based on GIM and the  
482 Bayesian Information Criterion to evaluate the spatiotemporal variability of net ecosystem  
483 exchange (NEE) as simulated by TBMs against atmospheric CO<sub>2</sub> concentration measurements  
484 from 35 towers in North America in 2008. We demonstrate the applicability of this approach by  
485 evaluating 4 prototypical TBMs participating in the North American Carbon Program Regional  
486 Interim Synthesis (NACP RIS).

487 We first design a series of synthetic data experiments in which the underlying fluxes are  
488 prescribed, to test the proposed approach and examine whether, when, and where atmospheric  
489 measurements are sensitive to, and hence can constrain, the spatiotemporal variability simulated  
490 by different TBMs. We find that due to the poor data coverage and weaker biospheric signals,  
491 current atmospheric observations cannot be used to reliably assess the flux spatiotemporal  
492 variability in the Tundra, Desert and Xeric Shrublands, and Tropical and Subtropical biomes.  
493 The remaining four biomes (i.e., Temperate Broadleaf and Mixed Forests, Temperate  
494 Grasslands, Savannas and Shrublands, Boreal Forests and Taiga, and Temperate Coniferous

495 Forest), however, are found to be well constrained by atmospheric data. Over these four biomes,  
496 the synthetic data experiments suggest that the proposed model selection approach, combined  
497 with the available atmospheric data, are able to identify the TBMs that represent a substantial  
498 portion of the underlying flux variability, as well as differentiate among multiple competing  
499 TBMs.

500 We further test and demonstrate the application of the approach by evaluating the performance of  
501 four prototypical TBMs that have been extensively assessed in literature using actual  
502 atmospheric observations. We find that conclusions about model performance are consistent with  
503 existing literature for cases where results are comparable, further supporting the applicability of  
504 our approach. Those results include that 1) TBMs represent fluxes best during the growing  
505 season (May-September) and least consistently with atmospheric observations during the  
506 transition seasons, especially in October and November; and that 2) TBMs appear to perform  
507 best over the Temperate Broadleaf and Mixed Forests biome. The experiments performed here  
508 also lead to new conclusions about the examined TBMs. For example, results show that SiB3  
509 and ORCHIDEE appear to represent the flux variability within individual biomes and months  
510 better relative to CASA-GFED and VEGAS2. In addition, this approach has the potential to link  
511 model performance with environmental processes, making it possible to test the hypothesis that  
512 seasonal differences in TBM performance reflect models' ability to represent the seasonal  
513 variability in the dominant environmental controls on fluxes.

514 The comparison conducted here only included four TBMs, and was intended primarily as a  
515 demonstration of the proposed approach. Furthermore, these four TBMs were not run using a  
516 uniform experimental protocol (Huntzinger et al., 2012), precluding any conclusive results about  
517 linkages between model performance and model structure. Applying the approach presented



518 here to a larger ensemble of models, ideally following a uniform simulation protocol, therefore  
519 represents a logical next step.

## 520 **Acknowledgments**

521 The authors thank the biospheric modelers participating the NACP Regional Interim Synthesis,  
522 and specifically Ning Zeng, Ian Baker, Nicolas Viovy, and James Randerson who provided the  
523 model results used in the analysis presented here. We thank Deborah Huntzinger for downscaling  
524 the CASA-GFED and VEGAS2 fluxes to 3-hourly temporal resolution. Atmospheric and  
525 Environmental Research (AER), and in particular Thomas Nehrkorn, John Henderson, and  
526 Janusz Eluszkiewicz, performed the WRF-STILT simulations and provided the sensitivity  
527 footprints. We acknowledge various data providers for the continuous in-situ CO<sub>2</sub> measurements,  
528 as well as Sharon Gourdji and Kim Mueller for their earlier efforts in collecting and processing  
529 those datasets. This work is funded by the National Aeronautics and Space Administration  
530 (NASA) under Grant No. NNX12AB90G and No. NNX12AM97G.

531

## 532 **Reference**

- 533 Baker, I., Denning, A. S., Hanan, N., Prihodko, L., Uliasz, M., Vidale, P.-L., Davis, K., and  
534 Bakwin, P.: Simulated and observed fluxes of sensible and latent heat and CO<sub>2</sub> at the  
535 WLEF-TV tower using SiB2.5, *Global Change Biology*, 9, 1262-1277, 2003.
- 536 Baker, I. T., Prihodko, L., Denning, A. S., Goulden, M., Miller, S., and Da Rocha, H. R.:  
537 Seasonal drought stress in the Amazon: Reconciling models and observations, *Journal of*  
538 *Geophysical Research: Biogeosciences*, doi: 10.1029/2007JG000644, 2008..
- 539 Balzarolo, M., Boussetta, S., Balsamo, G., Beljaars, A., Maignan, F., Calvet, J. C., Lafont, S.,  
540 Barbu, A., Poulter, B., Chevallier, F., Szczypta, C., and Papale, D.: Evaluating the  
541 potential of large-scale simulations to predict carbon fluxes of terrestrial ecosystems over  
542 a European Eddy Covariance network, *Biogeosciences*, 11, 2661-2678, 2014.
- 543 Beer, C., Reichstein, M., Tomelleri, E., Ciais, P., Jung, M., Carvalhais, N., Rödenbeck, C.,  
544 Arain, M. A., Baldocchi, D., Bonan, G. B., Bondeau, A., Cescatti, A., Lasslop, G.,  
545 Lindroth, A., Lomas, M., Luyssaert, S., Margolis, H., Oleson, K. W., Roupsard, O.,  
546 Veenendaal, E., Viovy, N., Williams, C., Woodward, F. I., and Papale, D.: Terrestrial  
547 Gross Carbon Dioxide Uptake: Global Distribution and Covariation with Climate,  
548 *Science*, 329, 834-838, 2010.

549 Belshe, E. F., Schuur, E. a. G., and Bolker, B. M.: Tundra ecosystems observed to be CO<sub>2</sub>  
550 sources due to differential amplification of the carbon cycle, *Ecology Letters*, 16, 1307-  
551 1315, 2013.

552 Broquet, G., Chevallier, F., Bréon, F. M., Kadygrov, N., Alemanno, M., Apadula, F., Hammer,  
553 S., Haszpra, L., Meinhardt, F., Morguá, J. A., Necki, J., Piacentino, S., Ramonet, M.,  
554 Schmidt, M., Thompson, R. L., Vermeulen, A. T., Yver, C., and Ciais, P.: Regional  
555 inversion of CO<sub>2</sub> ecosystem fluxes from atmospheric measurements: reliability of the  
556 uncertainty estimates, *Atmos. Chem. Phys.*, 13, 9039-9056, 2013.

557 Canadell, J. G., Ciais, P., Gurney, K., Le Quéré, C., Piao, S., Raupach, M. R., and Sabine, C. L.:  
558 An International Effort to Quantify Regional Carbon Fluxes, *Eos, Transactions American*  
559 *Geophysical Union*, 92, 81-82, 2011.

560 Chatterjee, A., Michalak, A. M., Anderson, J. L., Mueller, K. L., and Yadav, V.: Toward reliable  
561 ensemble Kalman filter estimates of CO<sub>2</sub> fluxes, *Journal of Geophysical Research:*  
562 *Atmospheres*, 117, D22306, 2012.

563 Globalview-Co2: Cooperative Atmospheric Data Integration Project – Carbon Dioxide (2010),  
564 CD-ROM, NOAA ESRL, Boulder, Colorado, also available at: <ftp.cmdl.noaa.gov>, Path:  
565 [ccg/co2/GLOBALVIEW](ftp://ftp.cmdl.noaa.gov/ccg/co2/GLOBALVIEW/)]. NOAA Global Monitoring Division: Boulder, Colorado,  
566 U.S.A., 2010.

567 Göckede, M., Michalak, A. M., Vickers, D., Turner, D. P., and Law, B. E.: Atmospheric inverse  
568 modeling to constrain regional-scale CO<sub>2</sub> budgets at high spatial and temporal resolution,  
569 *Journal of Geophysical Research: Atmospheres*, 115, D15113, 2010.

570 Gourdji, S. M., Hirsch, A. I., Mueller, K. L., Yadav, V., Andrews, A. E., and Michalak, A. M.:  
571 Regional-scale geostatistical inverse modeling of North American CO<sub>2</sub> fluxes: a  
572 synthetic data study, *Atmos. Chem. Phys.*, 10, 6151-6167, 2010.

573 Gourdji, S. M., Mueller, K. L., Yadav, V., Huntzinger, D. N., Andrews, A. E., Trudeau, M.,  
574 Petron, G., Nehrkorn, T., Eluszkiewicz, J., Henderson, J., Wen, D., Lin, J., Fischer, M.,  
575 Sweeney, C., and Michalak, A. M.: North American CO<sub>2</sub> exchange: inter-comparison of  
576 modeled estimates with results from a fine-scale atmospheric inversion, *Biogeosciences*,  
577 9, 457-475, 2012.

578 Hayes, D. J., Turner, D. P., Stinson, G., Mcguire, A. D., Wei, Y., West, T. O., Heath, L. S., De  
579 Jong, B., Mcconkey, B. G., Birdsey, R. A., Kurz, W. A., Jacobson, A. R., Huntzinger, D.  
580 N., Pan, Y., Post, W. M., and Cook, R. B.: Reconciling estimates of the contemporary  
581 North American carbon balance among terrestrial biosphere models, atmospheric  
582 inversions, and a new approach for estimating net ecosystem exchange from inventory-  
583 based data, *Global Change Biology*, 18, 1282-1299, 2012.

584 Hoeting, J. A., Davis, R. A., Merton, A. A., and Thompson, S. E.: Model Selection For  
585 Geostatistical Models, *Ecological Applications*, 16, 87-98, 2006.

586 Huntzinger, D. N., Gourdji, S. M., Mueller, K. L., and Michalak, A. M.: The utility of  
587 continuous atmospheric measurements for identifying biospheric CO<sub>2</sub> flux variability,  
588 *Journal of Geophysical Research: Atmospheres*, 116, D06110, 2011.

589 Huntzinger, D. N., Post, W. M., Wei, Y., Michalak, A. M., West, T. O., Jacobson, A. R., Baker,  
590 I. T., Chen, J. M., Davis, K. J., Hayes, D. J., Hoffman, F. M., Jain, A. K., Liu, S.,  
591 Mcguire, A. D., Neilson, R. P., Potter, C., Poulter, B., Price, D., Raczka, B. M., Tian, H.  
592 Q., Thornton, P., Tomelleri, E., Viovy, N., Xiao, J., Yuan, W., Zeng, N., Zhao, M., and  
593 Cook, R.: North American Carbon Program (NACP) regional interim synthesis:  
594 Terrestrial biospheric model intercomparison, *Ecological Modelling*, 232, 144-157, 2012.

595 Keenan, T. F., Baker, I., Barr, A., Ciais, P., Davis, K., Dietze, M., Dragoni, D., Gough, C. M.,  
596 Grant, R., Hollinger, D., Hufkens, K., Poulter, B., Mccaughey, H., Raczka, B., Ryu, Y.,  
597 Schaefer, K., Tian, H., Verbeeck, H., Zhao, M., and Richardson, A. D.: Terrestrial  
598 biosphere model performance for inter-annual variability of land-atmosphere CO<sub>2</sub>  
599 exchange, *Global Change Biology*, 18, 1971-1987, 2012.

600 Kort, E. A., Eluszkiewicz, J., Stephens, B. B., Miller, J. B., Gerbig, C., Nehrkorn, T., Daube, B.  
601 C., Kaplan, J. O., Houweling, S., and Wofsy, S. C.: Emissions of CH<sub>4</sub> and N<sub>2</sub>O over the  
602 United States and Canada based on a receptor-oriented modeling framework and  
603 COBRA-NA atmospheric observations, *Geophysical Research Letters*, 35, L18808, 2008.

604 Krinner, G., Viovy, N., De Noblet-Ducoudré, N., Ogée, J., Polcher, J., Friedlingstein, P., Ciais,  
605 P., Sitch, S., and Prentice, I. C.: A dynamic global vegetation model for studies of the  
606 coupled atmosphere-biosphere system, *Global Biogeochemical Cycles*, 19, GB1015,  
607 2005.

608 Lin, J. C., Gerbig, C., Wofsy, S. C., Andrews, A. E., Daube, B. C., Davis, K. J., and Grainger, C.  
609 A.: A near-field tool for simulating the upstream influence of atmospheric observations:  
610 The Stochastic Time-Inverted Lagrangian Transport (STILT) model, *Journal of*  
611 *Geophysical Research: Atmospheres*, 108, 4493, 2003.

612 Mcguire, A. D., Christensen, T. R., Hayes, D., Heroult, A., Euskirchen, E., Kimball, J. S.,  
613 Koven, C., Lafleur, P., Miller, P. A., Oechel, W., Peylin, P., Williams, M., and Yi, Y.: An  
614 assessment of the carbon balance of Arctic tundra: comparisons among observations,  
615 process models, and atmospheric inversions, *Biogeosciences*, 9, 3185-3204, 2012.

616 Mckain, K., Wofsy, S. C., Nehrkorn, T., Eluszkiewicz, J., Ehleringer, J. R., and Stephens, B. B.:  
617 Assessment of ground-based atmospheric observations for verification of greenhouse gas  
618 emissions from an urban region, *Proceedings of the National Academy of Sciences*, 109,  
619 8423-8428, 2012.

620 Michalak, A. M., Bruhwiler, L., and Tans, P. P.: A geostatistical approach to surface flux  
621 estimation of atmospheric trace gases, *Journal of Geophysical Research: Atmospheres*,  
622 109, D14109, 2004.

623 Miles, N. L., Richardson, S. J., Davis, K. J., Lauvaux, T., Andrews, A. E., West, T. O., Bandaru,  
624 V., and Crosson, E. R.: Large amplitude spatial and temporal gradients in atmospheric  
625 boundary layer CO<sub>2</sub> mole fractions detected with a tower-based network in the U.S. upper  
626 Midwest, *Journal of Geophysical Research: Biogeosciences*, 117, G01019, 2012.

627 Mueller, K. L., Yadav, V., Curtis, P. S., Vogel, C., and Michalak, A. M.: Attributing the  
628 variability of eddy-covariance CO<sub>2</sub> flux measurements across temporal scales using  
629 geostatistical regression for a mixed northern hardwood forest, *Global Biogeochemical*  
630 *Cycles*, 24, 2010.

631 Ogle, S., Davis, K., Andrews, A., West, T., Cook, R., Parkin, R., Morisette, J., Verma, S., and  
632 Wofsy, S.: Science Plan: Mid-Continent Intensive Campaign. Greenbelt, Md., 2006.

633 Peters, W., Jacobson, A. R., Sweeney, C., Andrews, A. E., Conway, T. J., Masarie, K., Miller, J.  
634 B., Bruhwiler, L. M. P., Pétron, G., Hirsch, A. I., Worthy, D. E. J., Van Der Werf, G. R.,  
635 Randerson, J. T., Wennberg, P. O., Krol, M. C., and Tans, P. P.: An atmospheric  
636 perspective on North American carbon dioxide exchange: CarbonTracker, *Proceedings of*  
637 *the National Academy of Sciences*, 104, 18925-18930, 2007.

638 Pillai, D., Gerbig, C., Kretschmer, R., Beck, V., Karstens, U., Neininger, B., and Heimann, M.:  
639 Comparing Lagrangian and Eulerian models for CO<sub>2</sub> transport – a step towards Bayesian  
640 inverse modeling using WRF/STILT-VPRM, *Atmos. Chem. Phys.*, 12, 8979-8991, 2012.

641 Ping, C.-L., Michaelson, G. J., Jorgenson, M. T., Kimble, J. M., Epstein, H., Romanovsky, V. E.,  
642 and Walker, D. A.: High stocks of soil organic carbon in the North American Arctic  
643 region, *Nature Geosci*, 1, 615-619, 2008.

644 Raczka, B. M., Davis, K. J., Huntzinger, D. N., Neilson, R., Poulter, B., Richardson, A., Xiao, J.,  
645 Baker, I., Ciais, P., Keenan, T. F., Law, B., Post, W. M., Ricciuto, D., Schaefer, K., Tian,  
646 H., Tomelleri, E., Verbeeck, H., and Viovy, N.: Evaluation of continental carbon cycle  
647 simulations with North American flux tower observations, *Ecological Monographs*, 83,  
648 531-556, 2013.

649 Richardson, A. D., Anderson, R. S., Arain, M. A., Barr, A. G., Bohrer, G., Chen, G., Chen, J. M.,  
650 Ciais, P., Davis, K. J., Desai, A. R., Dietze, M. C., Dragoni, D., Garrity, S. R., Gough, C.  
651 M., Grant, R., Hollinger, D. Y., Margolis, H. A., McCaughey, H., Migliavacca, M.,  
652 Monson, R. K., Munger, J. W., Poulter, B., Raczka, B. M., Ricciuto, D. M., Sahoo, A. K.,  
653 Schaefer, K., Tian, H., Vargas, R., Verbeeck, H., Xiao, J., and Xue, Y.: Terrestrial  
654 biosphere models need better representation of vegetation phenology: results from the  
655 North American Carbon Program Site Synthesis, *Global Change Biology*, 18, 566-584,  
656 2012.

657 Sasai, T., Ichii, K., Yamaguchi, Y., and Nemani, R.: Simulating terrestrial carbon fluxes using  
658 the new biosphere model “biosphere model integrating eco-physiological and  
659 mechanistic approaches using satellite data” (BEAMS), *Journal of Geophysical Research:*  
660 *Biogeosciences*, 110, G02014, 2005.

661 Schaefer, K., Schwalm, C. R., Williams, C., Arain, M. A., Barr, A., Chen, J. M., Davis, K. J.,  
662 Dimitrov, D., Hilton, T. W., Hollinger, D. Y., Humphreys, E., Poulter, B., Raczka, B. M.,  
663 Richardson, A. D., Sahoo, A., Thornton, P., Vargas, R., Verbeeck, H., Anderson, R.,  
664 Baker, I., Black, T. A., Bolstad, P., Chen, J., Curtis, P. S., Desai, A. R., Dietze, M.,  
665 Dragoni, D., Gough, C., Grant, R. F., Gu, L., Jain, A., Kucharik, C., Law, B., Liu, S.,  
666 Lokipitiya, E., Margolis, H. A., Matamala, R., McCaughey, J. H., Monson, R., Munger, J.  
667 W., Oechel, W., Peng, C., Price, D. T., Ricciuto, D., Riley, W. J., Roulet, N., Tian, H.,  
668 Tonitto, C., Torn, M., Weng, E., and Zhou, X.: A model-data comparison of gross  
669 primary productivity: Results from the North American Carbon Program site synthesis,  
670 *Journal of Geophysical Research: Biogeosciences*, 117, G03010, 2012.

671 Schuur, E. a. G., Vogel, J. G., Crummer, K. G., Lee, H., Sickman, J. O., and Osterkamp, T. E.:  
672 The effect of permafrost thaw on old carbon release and net carbon exchange from  
673 tundra, *Nature*, 459, 556-559, 2009.

674 Schwalm, C. R., Williams, C. A., Schaefer, K., Anderson, R., Arain, M. A., Baker, I., Barr, A.,  
675 Black, T. A., Chen, G., Chen, J. M., Ciais, P., Davis, K. J., Desai, A., Dietze, M.,  
676 Dragoni, D., Fischer, M. L., Flanagan, L. B., Grant, R., Gu, L., Hollinger, D., Izaurrealde,  
677 R. C., Kucharik, C., Lafleur, P., Law, B. E., Li, L., Li, Z., Liu, S., Lokupitiya, E., Luo,  
678 Y., Ma, S., Margolis, H., Matamala, R., McCaughey, H., Monson, R. K., Oechel, W. C.,  
679 Peng, C., Poulter, B., Price, D. T., Ricciuto, D. M., Riley, W., Sahoo, A. K., Sprintsin, M.,  
680 Sun, J., Tian, H., Tonitto, C., Verbeeck, H., and Verma, S. B.: A model-data  
681 intercomparison of CO<sub>2</sub> exchange across North America: Results from the North  
682 American Carbon Program site synthesis, *Journal of Geophysical Research:*  
683 *Biogeosciences*, 115, 566-584, 2010.

684 Shiga, Y. P., Michalak, A. M., Gourdji, S. M., Mueller, K. L., and Yadav, V.: Detecting fossil  
685 fuel emissions patterns from subcontinental regions using North American in situ CO<sub>2</sub>  
686 measurements, *Geophysical Research Letters*, 41, 2014GL059684, 2014.

687 Skamarock, W. C. and Klemp, J. B.: A time-split nonhydrostatic atmospheric model for weather  
688 research and forecasting applications, *J. Comput. Phys.*, 227, 3465-3485, 2008.

689 Tarnocai, C., Canadell, J. G., Schuur, E. a. G., Kuhry, P., Mazhitova, G., and Zimov, S.: Soil  
690 organic carbon pools in the northern circumpolar permafrost region, *Global*  
691 *Biogeochemical Cycles*, 23, GB2023, 2009.

692 Turner, D. P., Göckede, M., Law, B. E., Ritts, W. D., Cohen, W. B., Yang, Z., Hudiburg, T.,  
693 Kennedy, R., and Duane, M.: Multiple constraint analysis of regional land–surface  
694 carbon flux, *Tellus B*, 63, 207-221, 2011.

695 Van Der Werf, G. R., Randerson, J. T., Giglio, L., Collatz, G. J., Kasibhatla, P. S., and Arellano  
696 Jr, A. F.: Interannual variability in global biomass burning emissions from 1997 to 2004,  
697 *Atmos. Chem. Phys.*, 6, 3423-3441, 2006.

698 Yadav, V., Mueller, K. L., Dragoni, D., and Michalak, A. M.: A geostatistical synthesis study of  
699 factors affecting gross primary productivity in various ecosystems of North America,  
700 *Biogeosciences*, 7, 2655-2671, 2010.

701 Yadav, V., Mueller, K. L., and Michalak, A. M.: A backward elimination discrete optimization  
702 algorithm for model selection in spatio-temporal regression models, *Environmental*  
703 *Modelling & Software*, 42, 88-98, 2013.

704 Zeng, N., Mariotti, A., and Wetzel, P.: Terrestrial mechanisms of interannual CO<sub>2</sub> variability,  
705 *Global Biogeochemical Cycles*, 19, GB1016, 2005.

706

707

708 **Figures**

709 Figure 1. North American biomes, modified from Olson (2001), as defined for the case studies;  
710 green triangles indicate the locations of atmospheric CO<sub>2</sub> measurement towers used in the  
711 analysis.

712 Figure 2. Illustration of the 1°×1° and 3-hourly spatiotemporal variability of NEE simulated by  
713 CASA-GFED for Boreal Forests and Taiga in July. A vector including these 1°×1° and 3-hourly  
714 fluxes corresponds to one ancillary variable (i.e. one column) in **X**)

715 Figure 3. Illustration of Synthetic Data (SD) case studies as described in Section 4.

716 Figure 4. Average numbers of months within each season for which the candidate TBM is  
717 selected for the SD-one-ØØ, SD-one-ØΕ and SD-one-ΞΕ case studies (Figure 3). Grey shading  
718 in SD-one-ΞΕ represents biomes that were determined not to be well constrained by available  
719 atmospheric data. DJF: December, January, February; MAM: March, April, May; JJA: June,  
720 July, August; SON: September, October, November. The criteria for grey areas includes: 1) no  
721 models are selected in one season; or 2) the overall model selection is less than 50% in a year.

722 Figure 5. Average numbers of months within each season for which the candidate TBM is  
723 selected for the SD-all-ΞΕ case study (Figure 3). Grey shading represents biomes that were  
724 determined not to be well constrained by available atmospheric data. DJF: December, January,  
725 February; MAM: March, April, May; JJA: June, July, August; SON: September, October,  
726 November.

727 Figure 6. The correlation coefficient of NEE spatiotemporal series as simulated by different  
728 TBMs throughout 2008 for the four biomes better constrained by available atmospheric  
729 observations. TGSS: Temperate Grasslands, Savannas, Shrublands; Bore: Boreal Forests and  
730 Taiga; TCoF: Temperate Coniferous Forests; TBMF: Temperate Broadleaf and Mixed Forests.

731 Figure 7. Number of TBMs that are selected for each biome-month in the RD-one-ΞΕ cases  
732 study. Grey shading represents biomes that were determined not to be well constrained by  
733 available atmospheric data.

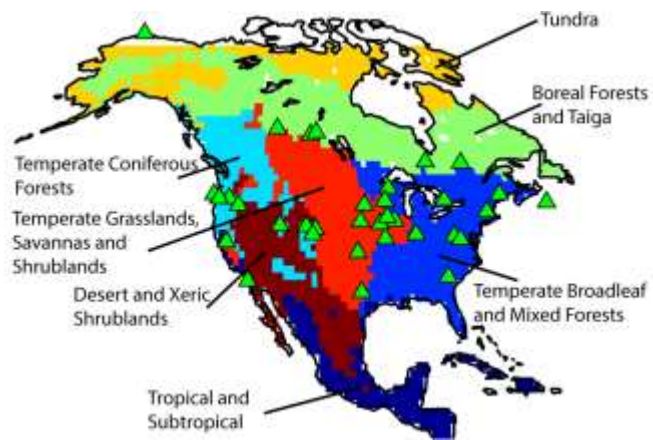
734 Figure 8. The TBM that explains the most variability in atmospheric measurements for a given  
735 biome-month, as identified by the RD-all-ΞΕ experiment. Grey shading represents biomes that  
736 were determined not to be well constrained by available atmospheric data.

737

738

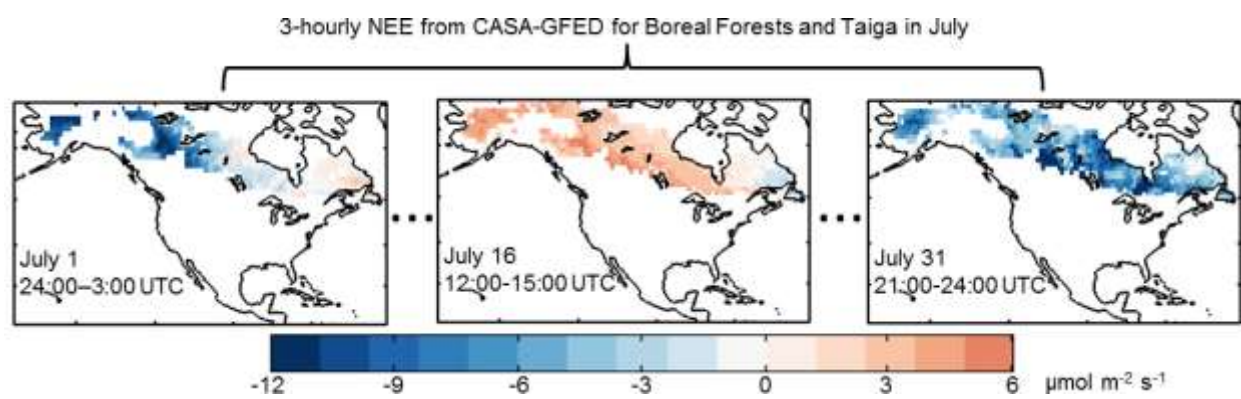
Figures

739 Figure 1



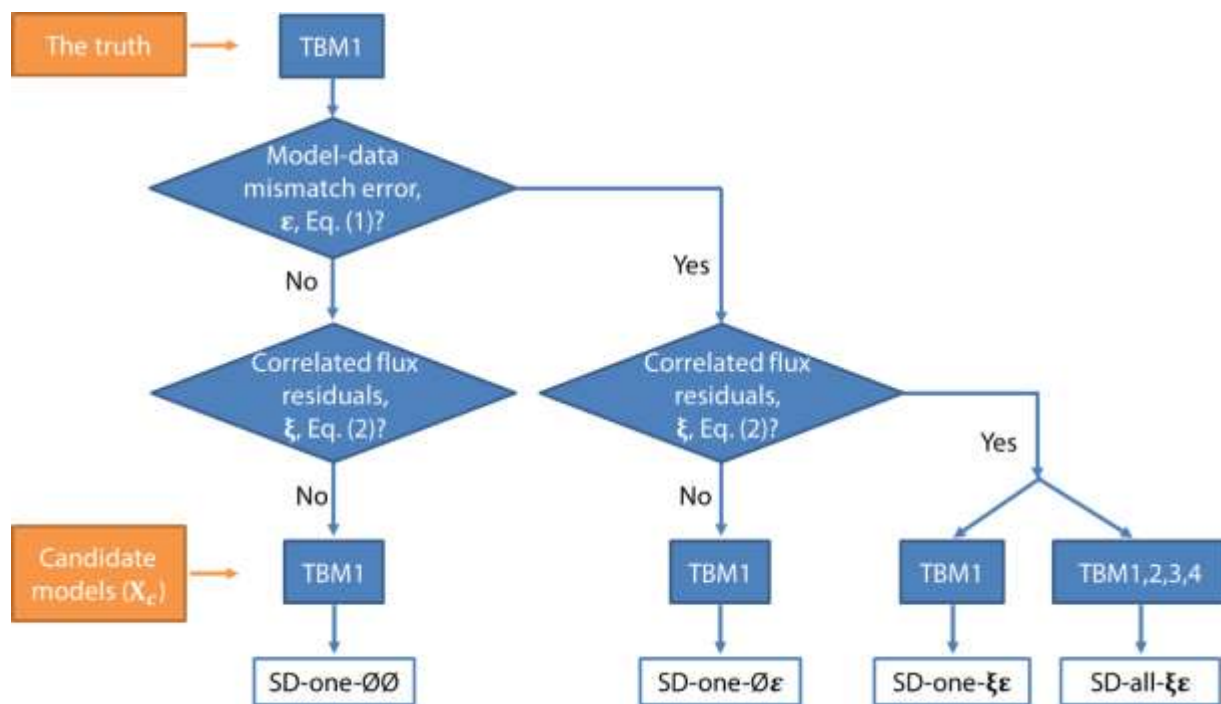
740

741 Figure 2



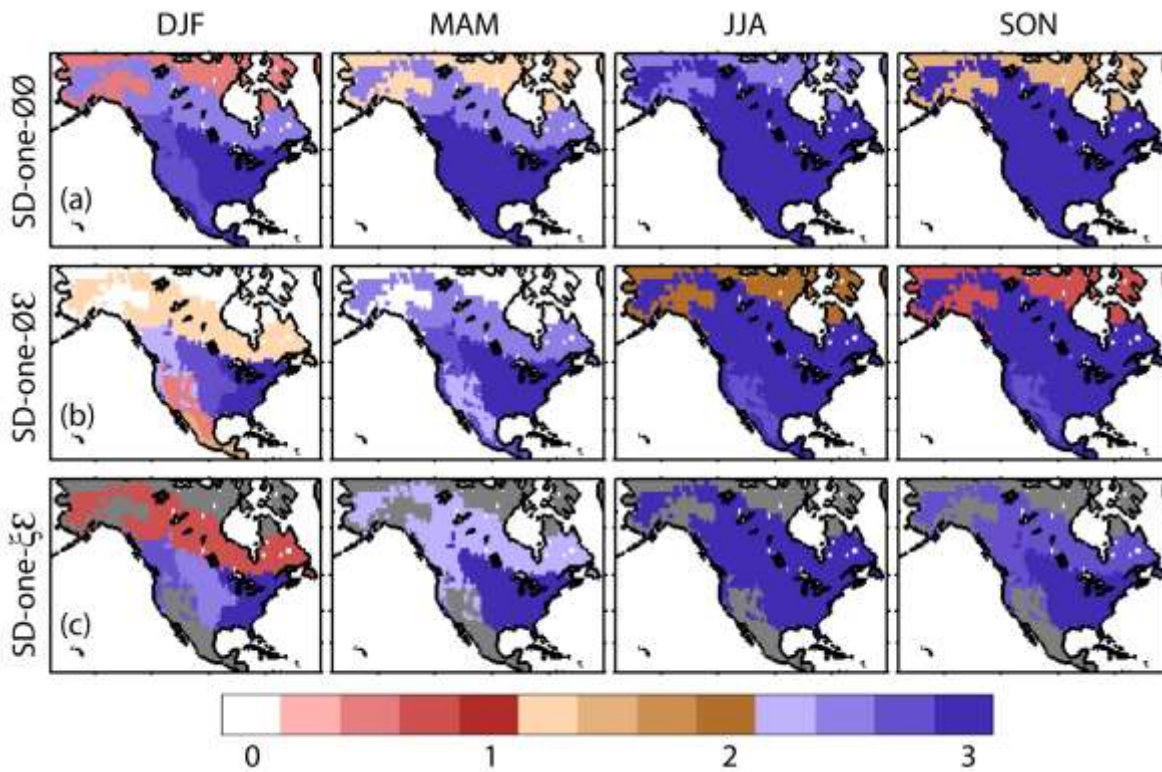
742

743 Figure 3



744

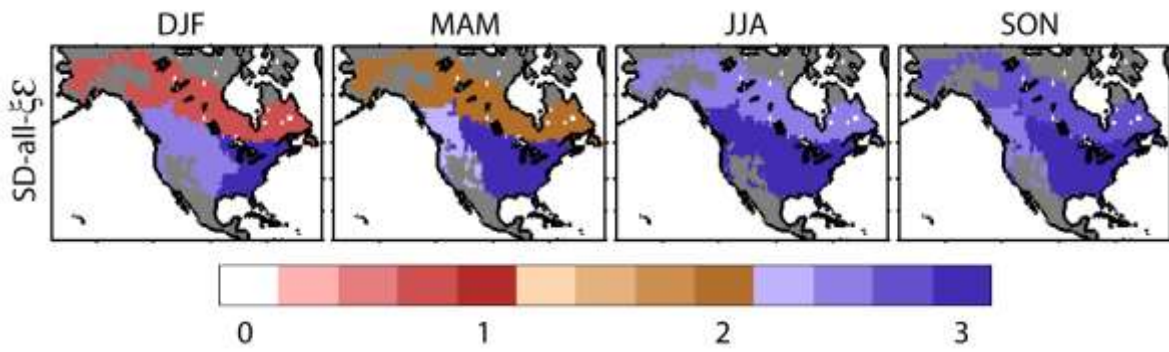
745 Figure 4



746

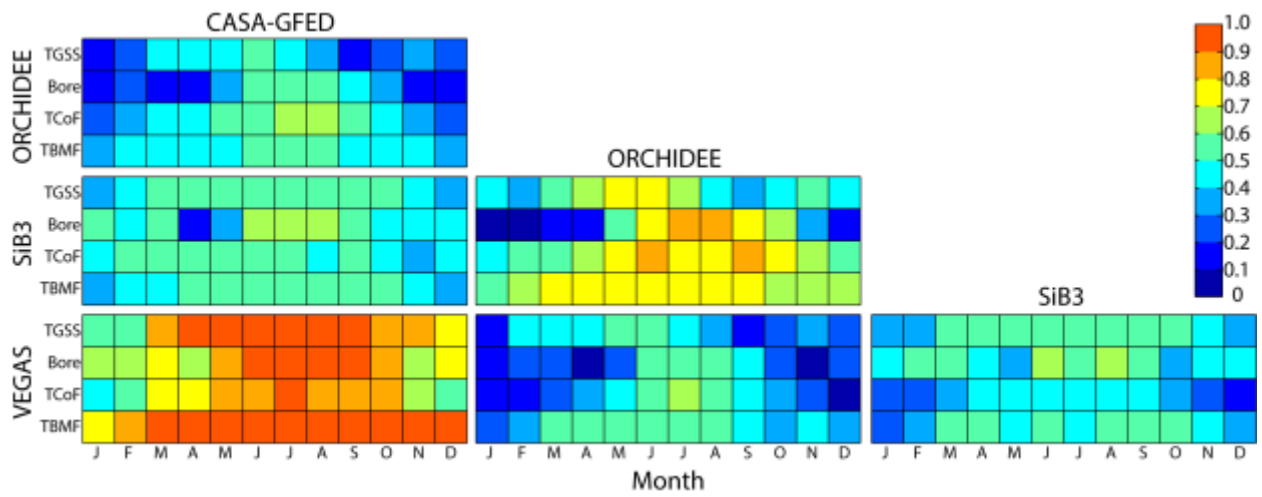


747 Figure 5



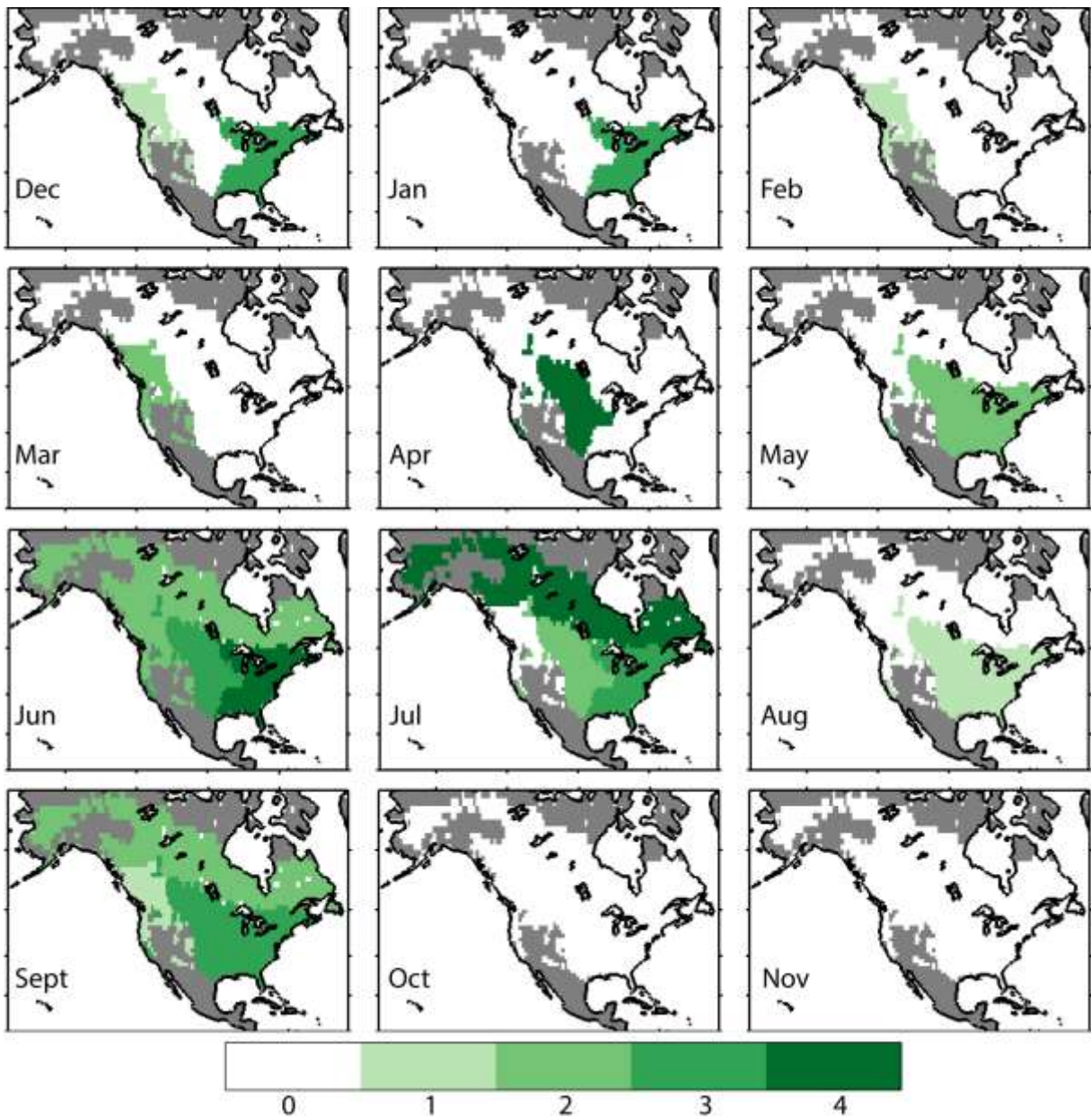
748

749 Figure 6



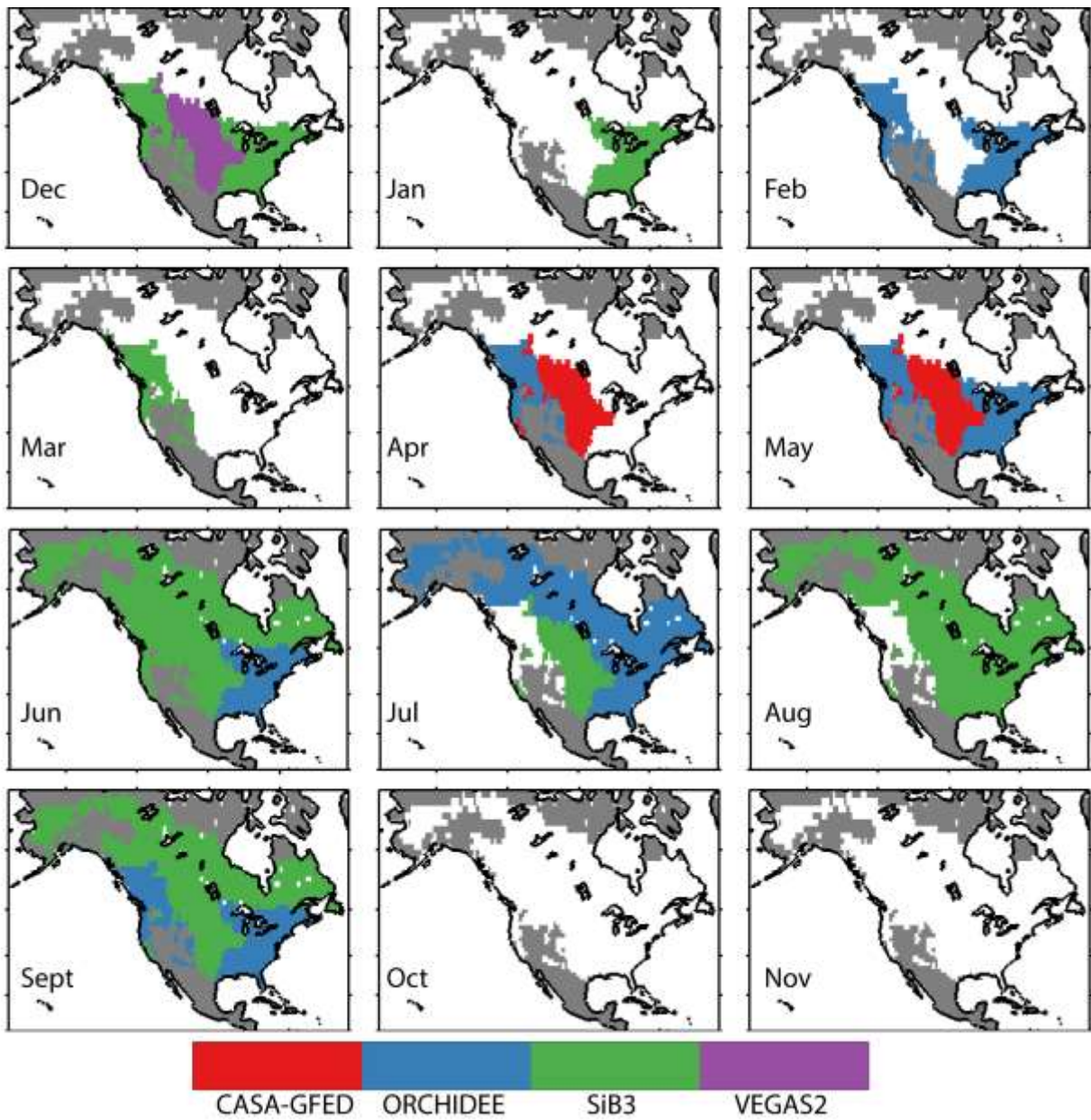
750

751 Figure 7



752

753 Figure 8



754