

1 **“Variable localization” in an Ensemble Kalman Filter:**
2 **application to the carbon cycle data assimilation**

3
4 ¹Ji-Sun Kang (jskang@atmos.umd.edu),

5 ¹Eugenia Kalnay(ekalnay@atmos.umd.edu),

6 ²Junjie Liu(jjliu@atmos.berkeley.edu),

7 ²Inez Fung (ifung@berkeley.edu),

8 ¹Takemasa Miyoshi (miyoshi@atmos.umd.edu), and

9 ¹Kayo Ide (ide@atmos.umd.edu)

10 ¹Department of Atmospheric and Oceanic Science, University of Maryland, College Park,

11 Maryland, 20742

12 ²Department of Earth and Planetary Science, University of California, Berkeley, 94720

13 **Abstract**

14 In Ensemble Kalman Filter (EnKF), space localization is used to reduce the impact
15 of long-distance sampling errors in the ensemble estimation of the forecast error
16 covariance. When two variables are not physically correlated, their error covariance is
17 still estimated by the ensemble, and therefore it is dominated by sampling errors. We
18 introduce a “variable localization” method, zeroing out such covariances between unrelated
19 variables to the problem of assimilating carbon dioxide concentrations into a dynamical
20 model using the Local Ensemble Transform Kalman Filter (LETKF) in an Observing
21 System Simulation Experiments (OSSE) framework. A system where meteorological and
22 carbon variables are simultaneously assimilated is used to estimate surface carbon fluxes
23 that are not directly observed. A range of covariance structures are explored for the
24 LETKF, with emphasis on configurations allowing non-zero error covariance between
25 carbon variables and the wind field, which affects transport of atmospheric CO₂, but not
26 between CO₂ and the other meteorological variables. Such “variable localization” scheme
27 zeroes out the background error covariance among prognostic variables that are not
28 physically related, thus reducing sampling errors. Results from the identical twin
29 experiments show that the performance in the estimation of surface carbon fluxes obtained

30 using “variable localization” is much better than that using a standard full covariance
31 approach. The relative improvement increases when the surface fluxes change with time
32 and model error becomes significant.

33 **1. Introduction**

34 An Observing System Simulation Experiments (OSSEs) system for carbon cycle data
35 assimilation has been created in parallel to a similar system that uses real meteorological
36 and CO₂ observations, and a state of the art model (Kang, 2009; Liu et al., CO₂ transport
37 uncertainties from the uncertainties in meteorological fields, submitted to *Geophys. Res.*
38 *Lett.*, 2011). The ultimate goal of these parallel projects is to estimate not only
39 atmospheric CO₂ but also surface carbon fluxes. This is a challenging problem plagued
40 with obstacles whose origin frequently cannot be even identified using real data and
41 without knowing the “truth”. In the course of performing OSSEs, we have found several
42 algorithms that can substantially improve the results. The focus of this paper is on one of
43 these algorithms, “variable localization” that reduces sampling errors and can be also
44 applied to other problems in data assimilation.

45 The Local Ensemble Transform Kalman Filter [LETKF, *Hunt et al.*, 2007], like other
46 Ensemble Kalman Filter (EnKF) methods [*Evensen*, 1994; *Houtekamer and Mitchell*, 2001;
47 *Anderson*, 2001; *Bishop et al.*, 2001; *Whitaker and Hamill*, 2002; *Ott et al.*, 2004, *Zupanski*,
48 2005, and others], produces an analysis using a multivariate background error covariance
49 matrix that contains an estimation of the error correlation between the dynamic variables.

50 When the variables are physically related to each other, the multivariate background error
51 estimation helps the analysis to efficiently correct the forecast errors. Indeed, *Liu et al.*
52 [2009] have shown that a multivariate assimilation of AIRS (Atmospheric Infrared
53 Sounder) humidity retrievals has lower wind analysis errors than the standard univariate
54 assimilation used in operational numerical weather prediction (NWP) system that does not
55 account for the error covariance between humidity and winds. However, standard
56 multivariate EnKF also allows for error covariances among model variables even if some of
57 those variables are not physically related to each other. In this case, the estimate of the
58 error covariance will be solely due to sampling errors.

59 The direct solution to reduce sampling errors would be to increase the number of
60 ensemble members, but this is not a practical solution because of computational and storage
61 requirements. It is common practice in EnKF with a limited ensemble size to introduce
62 “space localization” into the background error covariance [*Houtekamer and Mitchell, 2001*;
63 *Hamill et al., 2001*]. The background ensemble perturbations have error covariances that
64 are good estimates of real covariances over relatively short distances of up to about 500-
65 1000 km in the global NWP applications. At longer distances, the background errors are
66 still apparently correlated, but these correlations become dominated by sampling errors, and

67 can seriously harm the analysis. In the widely adopted technique of “space localization”
68 to solve the problem of long-distance spurious correlations, the background error
69 covariance terms are multiplied by an approximation of a Gaussian function that decreases
70 with the distance between the two grid points whose error covariance is being computed
71 and becomes negligible at distances greater than about 1000 km [*Gaspari and Cohn, 1999*].

72 In our carbon cycle data assimilation OSSEs, we apply a similar concept whenever two
73 variables are not *physically* related and therefore estimates of their error covariances are
74 spurious. In that case, we avoid spurious correlations by zeroing out these covariances
75 due to sampling errors. For example, atmospheric CO₂ concentration is determined by the
76 wind transport as well as by CO₂ surface fluxes. However, the evolution of the carbon
77 variables is likely to have a much less significant dependence on some other variables such
78 as the specific humidity or surface pressure. If this is the case, we zero out the error
79 covariance between the atmospheric CO₂ and both of specific humidity and surface
80 pressure in the analysis. This new methodology is denoted as “variable localization” by
81 analogy with space localization, because the background error covariances between
82 variables that are not physically linked in a significant way are zeroed out.

83 There are several previous studies that estimate surface carbon fluxes via data

84 assimilation. *Baker et al.* [2006, 2008] have applied a four-dimensional variational (4D-
85 Var) method to their OSSEs while *Peters et al.* [2005, 2007] and *Feng et al.* [2008] have
86 used EnKFs. *Peters et al.* [2005, 2007] have assimilated observations from the ground-
87 based stations whereas *Feng et al.* [2008] have used simulated observations of satellite
88 data, Orbiting Carbon Observatory [OCO; *Crisp et al.*, 2004]. Furthermore, as a part of
89 Global and regional Earth-system (Atmosphere) Monitoring project [GEMS, *Hollingsworth*
90 *et al.*, 2008], a system with a two-step approach has been built for a carbon cycle data
91 assimilation system: the first is to assimilate satellite and in situ data to monitor the
92 atmospheric CO₂ within a 4D-Var [*Engelen et al.*, 2009] and the second is a variational flux
93 inversion system [*Chevallier et al.*, 2009a and b]. These studies have shown the
94 meaningful results in estimating surface CO₂ fluxes from the atmospheric CO₂
95 concentration observations of in-situ data as well as the satellite data. On the other hand,
96 *Zupanski et al.* [2007] have applied Maximum Likelihood Ensemble Filter [*Zupanski,*
97 2005] to bias estimation of surface CO₂ fluxes over a local area with several tower
98 observations.

99 In these studies, surface CO₂ fluxes are estimated by assimilating the observations of
100 atmospheric CO₂ concentration, but not by any direct observations of carbon fluxes or other

101 meteorological variables. In order to link the surface carbon fluxes with the atmospheric
102 CO₂ concentrations, these studies have used transport models that play an important role in
103 transferring information from the atmospheric CO₂ observations to corresponding changes
104 in the surface flux of carbon. Also, these studies need *a priori* information about the
105 carbon variables as an initial guess that is pre-calculated using independent observations or
106 model simulations because the problem of determining surface carbon fluxes is otherwise
107 ill-posed [Enting, 2002]. The surface CO₂ fluxes are determined by minimizing the
108 squared normalized difference between the simulated CO₂ concentration and the observed
109 CO₂, and *a priori* error for the atmospheric CO₂ concentrations and fluxes based on their
110 error covariances.

111 So far, data assimilation studies of carbon fluxes have not yet accounted for transport
112 errors in the atmospheric CO₂ forecast that can be caused by both the imperfections of the
113 transport model and the uncertainty of wind analysis which drives the transport model, even
114 though many studies [Gurney *et al.*, 2004; Baker *et al.*, 2008; Stephens *et al.*, 2007;
115 Miyazaki, 2009] have found that the accuracy of atmospheric CO₂ forecasts depends on
116 those transport errors. The Bayesian synthesis approaches, usually referred to as inversion
117 modeling, [Bousquet *et al.*, 2000; Gurney *et al.*, 2004; Rödenbeck *et al.*, 2003] have been

118 used to estimate surface carbon fluxes prior to the advent of the data assimilation studies
119 discussed above, and also have the same limitation stemming from unresolved transport
120 errors. Indeed, some studies [*Gurney et al.*, 2004; *Baker et al.*, 2008; *Stephens et al.*,
121 2007] have pointed out that the transport errors can cause biases in both the atmospheric
122 CO₂ analysis and the surface CO₂ flux estimation. Notably, *Miyazaki* [2009] shows a
123 significant contrast in the results of atmospheric CO₂ forecast obtained using wind fields of
124 different accuracies. This result strongly emphasizes the importance of wind uncertainty
125 in carbon cycle data assimilation.

126 As a complement of our real-data experiments (Liu et al., 2011), we present here a
127 similar OSSE carbon cycle data assimilation system that simultaneously assimilates the
128 observations of meteorological variables (wind, temperature, humidity, and surface
129 pressure) and atmospheric CO₂. The system analyzes not only these meteorological
130 variables and atmospheric CO₂, but also the surface CO₂ fluxes. Since our method
131 generates the analysis of meteorological variables and carbon variables simultaneously, we
132 do not need to run a transport model for carbon variables in addition to running a forecast
133 model for meteorological variables. Besides, results of our method do not depend on *a*
134 *priori information* on the initial condition of carbon variables which, like the

135 meteorological variables, spin-up and converge even if the ensembles are started from
136 smooth random fields [Zupanski *et al.*, 2006]. Since our challenging ultimate goal is to
137 estimate surface CO₂ fluxes from a simultaneous analysis of meteorological variables and
138 carbon variables, we have tested new techniques to improve the ability of EnKF to reach
139 this goal. Here we introduce one of these techniques, “variable localization” that can be
140 usefully applied in any EnKF system.

141 Previous data assimilation estimates of surface carbon fluxes [Baker *et al.*, 2006, 2008;
142 Peters *et al.*, 2007; Feng *et al.*, 2008; Chevallier *et al.*, 2009a] can be considered as
143 belonging to the case of “carbon-univariate” analyses where the atmospheric CO₂
144 concentrations and the surface CO₂ fluxes are updated by themselves, without including
145 error correlations between these carbon variables and meteorological variables. In this
146 study, various types of analyses including error covariances are introduced and compared,
147 ranging from a standard “fully multivariate” analysis to a “carbon-univariate” analysis
148 within the LETKF framework. Because this is the first test of a new methodology, this
149 work is limited to OSSEs (in a “twin experiment” approach) in which the observations are
150 sampled from a “nature run” assumed to be the true evolution of the system and assimilated
151 using the LETKF. For simplicity, we also assume that the model is perfect except in the

152 last experiments where the surface fluxes of carbon are varied in the nature run but not in
153 the forecast model, and we focus on the impact of various “variable localization”
154 techniques.

155 The paper is organized as follows. Section 2 provides a description of the model used
156 for this study and the various “variable localization” schemes tested within the LETKF data
157 assimilation framework. Section 3 describes the experimental design. Results are shown
158 in section 4, and we summarize and discuss our findings in section 5.

159

160 **2. Methodology**

161 *2.1. Model: SPEEDY-C*

162 The SPEEDY model [Molteni, 2003] is an atmospheric global, primitive equation
163 general circulation model (AGCM) with simplified physical parameterization schemes that
164 is computationally efficient, but it maintains the basic characteristics of a state-of-the-art
165 AGCM with complex physics. The version used for this study has triangular truncation
166 T30 with 7 vertical sigma levels, and has five dynamical variables including zonal (U) and
167 meridional wind (V) components, temperature (T), specific humidity (q), and surface
168 pressure (Ps).

169 To incorporate atmospheric carbon dioxide concentration (C) and surface flux of
170 carbon dioxide (CF), the SPEEDY model is extended to the “SPEEDY-C” model, which
171 contains these carbon-related variables.

$$172 \quad \frac{\partial(C)}{\partial t} + \mathfrak{Z}(C) = CF \quad (1)$$

173 Equation (1) shows how the tendency of atmospheric CO_2 is calculated in SPEEDY-C,
174 where $\mathfrak{Z}(C)$ represents the atmospheric 3-dimensional transport and mixing. In this
175 study, chemical processes affecting atmospheric carbon dioxide are ignored since CO_2 in
176 the atmosphere is essentially inert. Moreover, there is no feedback between the integrated
177 CO_2 and the radiative properties of the SPEEDY-C model. Surface flux of carbon (CF) on
178 the right-hand side of Equation (1) provides sources and sinks of CO_2 . Carbon flux on the
179 surface (CF) is converted into the atmospheric CO_2 concentrations added to the lowest
180 layer of the model. In reality, several types of forcings make up this flux: fossil fuel
181 emission, land surface fluxes due to vegetation and land use change, and ocean fluxes. In
182 this OSSE study, we test the ability of our data assimilation systems with a “variable
183 localization” to estimate surface CO_2 fluxes, so for simplicity we assume that the CF is due
184 only to constant fossil fuel emissions for most of experiments. However, in the last
185 experiment (shown in Figure 10), we do allow for variable fluxes associated with

186 vegetation and ocean in the nature run but not in the forecast model run. Thus, the
 187 SPEEDY-C has six prognostic variables (U, V, T, q, Ps, C), along with either a constant or a
 188 variable forcing (CF) which is not changed by the model.

189 *2.2. The LETKF*

190 The LETKF is an ensemble Kalman filter method where the background error
 191 covariance \mathbf{P}^b among the variables can be estimated as

$$192 \quad \mathbf{P}^b = \frac{1}{K-1} \mathbf{X}^b \mathbf{X}^{bT} \quad (2)$$

193 where \mathbf{X}^b is the matrix whose columns contain a departure of each ensemble forecast
 194 ($\mathbf{x}^{b(i)}$) from the ensemble mean ($\bar{\mathbf{x}}^b$): the i -th column of \mathbf{X}^b is $\mathbf{x}^{b(i)} - \bar{\mathbf{x}}^b$, $\{i=1,2,\dots,K\}$,
 195 K is the number of ensemble members and \mathbf{x} is a state vector of dynamic variables at the
 196 model grids. The evolution of \mathbf{P}^b , which contains the background error covariance
 197 among the dynamic variables, is accounted for in every analysis step so that temporally and
 198 spatially varying uncertainties in the background are considered when analyzing variables.

199 The first step of the analysis is to compute \mathbf{X}^b . Then, the observation operator h is
 200 applied to the ensemble forecast \mathbf{x}^b to transform the background from the model grid
 201 space to the observation space, $\mathbf{y}^{b(i)} = h(\mathbf{x}^{b(i)})$. Let $\mathbf{Y}^b = \mathbf{y}^{b(i)} - \bar{\mathbf{y}}^b$, $\{i=1,2,\dots,K\}$ be
 202 the background perturbations in the observation space. Then, the estimation of the

203 background is ready to be compared with observations in the same space.

204 In order to produce the analysis at every grid point, the LETKF assimilates only
205 observations within a certain distance from each grid point so that the following analysis
206 computations are performed locally. The analysis mean, $\bar{\mathbf{x}}_{(l)}^a$, is given by

$$207 \quad \bar{\mathbf{x}}_{(l)}^a = \bar{\mathbf{x}}_{(l)}^b + \mathbf{X}_{(l)}^b \bar{\mathbf{w}}_{(l)}, \quad (3)$$

208 where $\bar{\mathbf{w}}_{(l)}$ is the mean weighting vector calculated by

$$209 \quad \bar{\mathbf{w}}_{(l)} = \tilde{\mathbf{P}}_{(l)}^a (\mathbf{Y}_{(l)}^b)^T \mathbf{R}_{(l)}^{-1} (\mathbf{y}_{(l)}^o - \bar{\mathbf{y}}_{(l)}^b). \quad (4)$$

210 Here, $\tilde{\mathbf{P}}_{(l)}^a = [(\mathbf{Y}_{(l)}^b)^T \mathbf{R}_{(l)}^{-1} (\mathbf{Y}_{(l)}^b) + (K-1)\mathbf{I}/\rho]^{-1}$ is the analysis error covariance in the
211 ensemble space, \mathbf{R} is the observation error covariance matrix, \mathbf{y}^o is the observation
212 vector, and ρ is the inflation factor (see section 2.3.3 for details), and the subscript (l)
213 means a quantity defined on a local region centered at the analysis grid point l . Within a
214 local region, *space localization* is carried out by multiplying the inverse observation error
215 covariance matrix $\mathbf{R}_{(l)}^{-1}$ by a factor that decays from one to zero as the distance of the
216 observations from the analysis grid point increases [Miyoshi, 2005, Hunt et al., 2007,
217 Greybush et al., 2010].

218 The analysis increment, $\bar{\mathbf{x}}_{(l)}^a - \bar{\mathbf{x}}_{(l)}^b$ (Eqns. 3 and 4), is given by the background
219 perturbation matrix multiplied by the weight vector which is a function of the

220 innovation, $\mathbf{y}_{(t)}^o - \bar{\mathbf{y}}_{(t)}^b$, and error statistics of both background and observation. Thus, the
 221 analysis reflects observational information more than background information if the
 222 background error is greater than the observation errors, and vice versa. In addition, the
 223 ensemble perturbations of the analysis are determined by

$$224 \quad \mathbf{X}_{(t)}^a = \mathbf{X}_{(t)}^b [(K-1)\tilde{\mathbf{P}}_{(t)}^a]^{-\frac{1}{2}} \quad (5)$$

225 With (5) we obtain the estimation of analysis uncertainty in addition to the analysis mean.
 226 The global analysis ensemble $\mathbf{x}^{a(i)}$, $\{i=1,2,\dots,K\}$, is formed by gathering the values
 227 obtained for $\bar{\mathbf{x}}_{(t)}^a$ and $\mathbf{X}_{(t)}^a$ at all the analysis grid points. (see *Hunt et al.* [2007] for
 228 more details and discussion on LETKF.)

229 *2.3. LETKF application to the SPEEDY-C: variable localization*

230 *2.3.1. Motivation*

231 In order to estimate not only the model prognostic variables (U, V, T, q, Ps, C) but also
 232 the unknown surface fluxes field (CF), we use an augmented state vector \mathbf{x} consisting of
 233 (U, V, T, q, Ps, C, CF) at all model grid points, where CF, like Ps, is defined at the model
 234 surface grid points. This augmentation enables the LETKF to directly estimate the
 235 parameter like any other (unobserved) variable through the background error covariance
 236 with the observed variables [*Baek et al.*, 2006; *Annan et al.*, 2005].

237 More sophisticated schemes can be designed by taking into account that dynamical
238 interactions of the augmented variables are not homogeneous in the SPEEDY-C. As
239 shown in Equation (1), atmospheric CO₂ (C) is advected by (U, V) and forced by surface
240 carbon fluxes (CF) but has no direct interaction with (T, q, Ps). In contrast, none of the
241 meteorological variables (U, V, T, q, Ps) is dynamically affected by C or CF while CF is
242 not affected by any of the dynamical variables (U, V, T, q, Ps), at least within the SPEEDY
243 model formulation. When sampling the standard fully multivariate background error
244 covariances using a finite-size ensemble (Figure 1a), however, spurious correlations may
245 arise between the variables. This motivates us to develop analysis schemes by grouping
246 the variables based on the idea of the localization according to the “*dynamical distance*
247 *between the variables*”. This “variable localization” attempts to manage the correlations
248 between the model variable groups, like the conventional localization attempts to suppress
249 the spurious correlation based on the “*physical*” distance.

250 Various analysis methods are possible according to the method used to group variables.
251 For example, if one groups an analysis state vector of only carbon variables (C, CF) and
252 the other of meteorological variables (U, V, T, q, Ps) separately (Figure 1e), the analysis of
253 carbon is determined by only assimilating atmospheric CO₂ observations *univariately for*

254 *carbon* (Equations 3-5). In this case, the surface CO₂ fluxes are updated by the
255 multiplication of a background perturbation matrix of surface CO₂ fluxes and the weight
256 vector (Equation 4) as calculated from the forecast and the observations of atmospheric
257 CO₂ concentrations. If the analysis state vector is designed to include other
258 meteorological variables in addition to (C, CF), then the analysis can reflect the
259 background error covariance among all those variables in order to estimate surface CO₂
260 fluxes *multivariately* (e.g., Figure 1a). Such an approach implies that the analysis allows
261 error information to flow from carbon to the meteorological variables in the state vector
262 and vice versa.

263 2.3.2. *Different covariance structures for analyses: variable localization*

264 In this study, we introduce variable localization and test five analysis methods
265 characterized by the \mathbf{P}^b configurations based on the “dynamical” distance between
266 variables (Figure 1). The first method is the standard *fully multivariate* analysis
267 (hereafter referred as *mult*) in which the errors of all dynamic variables are coupled in the
268 background error covariance (Fig 1a). This scheme (used in present EnKF methods)
269 allows errors in all variables to be potentially correlated with one another. As a result,
270 the system gives more weight to the atmospheric CO₂ observations whenever any of the

271 dynamic variables have a larger uncertainty in the background field. On the other hand,
272 the uncertainty of the carbon variables can also change the weight vector $\bar{\mathbf{w}}_{(t)}$ (Equations
273 3 and 4), which is shared among all dynamic variables. Analysis uncertainty of all
274 variables (U, V, T, q, Ps, C, CF) is determined by Equation (5).

275 The second method is based on the notion that in our model the surface fluxes are only
276 physically related to the atmospheric CO₂ but not to other dynamical variables. That is,
277 the white areas of background error covariance matrix \mathbf{P}^b in Figure 1(b) contain sampling
278 errors rather than any useful error correlations or covariances. Thus, we zero out those
279 white areas, the covariances between CF and all variables except atmospheric CO₂, in order
280 to eliminate sampling errors in these correlations (*localized-multivariate* analysis: ***L-mult***,
281 Fig 1b). This scheme has two separate analyses, one for (U, V, T, q, Ps, C) (in grey) and
282 the other for CF (in black). Analysis of surface carbon fluxes assimilates only
283 atmospheric CO₂ observations but not the observations of meteorological variables for
284 computing $\bar{\mathbf{w}}_{(t)}$ and the uncertainty of CF. In contrast, the analysis of the dynamic
285 variables except CF assimilates all available observations of (U, V, T, q, Ps, C) for
286 computing the other $\bar{\mathbf{w}}_{(t)}$ and their analysis uncertainty. In other words, a set of
287 Equations (4)-(5) is computed separately having $\mathbf{X}_1=(U, V, T, q, Ps, C)$ and $\mathbf{X}_2=(C, CF)$ to

288 get each weight vector $\bar{\mathbf{w}}_{(t)}$ and analysis error covariance $\tilde{\mathbf{P}}_{(t)}^a$ for updating (U, V, T, q,
289 Ps, C) and CF respectively. Eliminating spurious correlations with carbon fluxes is
290 especially important since we do not start the analysis with any *a priori* knowledge of
291 carbon variables. Because CF is not constrained by any direct observations, it is very
292 possible for the CF to degrade the analysis of the other variables due to bad initial values at
293 the initial stage of the *mult* analysis. Thus, poor initial conditions for carbon may degrade
294 the analysis of other meteorological variables in the *mult* analysis whereas *L-mult* prevents
295 initial carbon from poorly influencing the analysis of all other dynamic variables.

296 The third method is the *1-way multivariate* analysis (*Iway*) based on the notion that
297 wind uncertainties should be able to provide useful information to update carbon variables,
298 whereas the sampling error in the carbon variable is assumed to be too large to provide a
299 positive impact to the wind assimilation (Fig. 1c). In the *Iway* scheme, the atmospheric
300 CO₂ concentrations and surface CO₂ fluxes are updated using an error covariance that
301 includes the wind fields, while the wind and other atmospheric variables such as
302 temperature, specific humidity, and surface pressure are updated separately and are not
303 affected by these two carbon variables (Fig 1c). This scheme was also found useful by *Liu*
304 *et al.* [2009] when assimilating AIRS moisture retrievals.

305 The fourth method is also based on the *Iway* system, but zeroing out the background
306 error covariance between surface carbon fluxes and wind fields. We refer to this as the
307 *localized-Iway multivariate* analysis (***L-Iway***, Fig 1d), as in the case of the ***L-mult*** scheme.
308 It is based on the idea that winds transport atmospheric CO₂ but not surface carbon fluxes,
309 and thus their errors should be uncorrelated. Here, the resulting analysis of
310 meteorological variables should be exactly same as in *Iway* (Fig. 1c). The comparison of
311 the ***L-Iway*** method with *Iway* provides a measure of the direct impact of wind
312 uncertainties on the estimation of surface CO₂ fluxes.

313 The last method considered is the *Carbon-univariate* analysis (***C-univ***). In this
314 method, atmospheric CO₂ concentration and surface CO₂ fluxes are updated only by these
315 two variables themselves, unaffected by other atmospheric variables (Fig 1e). The
316 forecasts of atmospheric CO₂ are still driven by the ensemble of wind fields. Although the
317 ensemble transport of CO₂ provides some information about wind uncertainties to the
318 background state of atmospheric CO₂ in ***C-univ***, the transport error term is not explicitly
319 used for the carbon analysis.

320 2.3.3. *Inflation of the background covariance*

321 In practice, the ensemble forecast tends to underestimate the uncertainty in its state

322 estimate because of limited ensemble size, model errors and nonlinearities. To
323 compensate for this underestimation, it is necessary to inflate the background covariance
324 (or the analysis covariance) during each data assimilation cycle. For the inflation factor,
325 multiplicative inflation has been applied in this work [Anderson and Anderson, 1999].
326 This is carried out by multiplying the background perturbation from the ensemble mean by
327 a factor larger than one (ρ). It is common to tune this inflation parameter manually;
328 however, such tuning is expensive, and becomes infeasible if the inflation factor is allowed
329 to depend on space and time, and/or the variable. Since we have found that the carbon
330 variables require quite different inflation factors compared to the inflation for the
331 meteorological variables, the adaptive inflation estimation introduced by *Li et al.* [2009]
332 has been used to estimate the inflation factors for the meteorological variables, on the one
333 hand, and the atmospheric CO₂ concentration on the other. *Li et al.* [2009] estimated
334 simultaneously the adaptive inflation and observation errors, using the equations derived by
335 *Desroziers et al.* [2005]. Here we assume that the observation error statistics are correct,
336 and we calculate the inflation adaptively for each vertical layer separately. Moreover, for
337 the atmospheric CO₂ in the lowest layer, we calculate and apply two separate inflation
338 factors over the land and the ocean areas. The methodology of *Li et al.* [2009] compares

339 the analysis increment (analysis minus background) and the observation increment
340 (background minus observation) with the expected values in observation space. Thus, that
341 methodology is only available for variables having observations, which means we need to
342 apply a different method for the inflation of surface carbon fluxes (CF). Our approach for
343 CF is similar to the covariance relaxation method of *Zhang et al.* [2004], except that we let
344 the analysis perturbations maintain the same spread as the background. More details
345 about the adaptive inflation methods can be found in *Li et al.* [2009] and *Zhang et al.*
346 [2004].

347

348 **3. Experimental design: Observing System Simulation Experiments (OSSEs)**

349 In the OSSEs, the SPEEDY-C model with a total constant fossil fuel emission of
350 6PgC/yr [*Andres et al.*, 1996; Figure 2a] is used to create the “nature run” assumed to be
351 the true state in this study (but we also perform an OSSE with varying surface fluxes,
352 obtained with a model with interactive vegetation, see Figure 10). Simulated observations
353 are then obtained from this “nature run” by adding random observational errors. Standard
354 deviations of the simulated observation errors are listed in Table 1. For the atmospheric
355 variables, the observations have the spatial distribution of the rawinsonde network, with

356 about 9% coverage of grid points globally (Figure 3a), with more observations in the
 357 Northern Hemisphere mid-latitudes.

358 Atmospheric CO₂ concentration is assumed to be observed from three different
 359 measurements: one comes from 18 *in situ* data locations which have continuous records of
 360 CO₂ concentration near the surface (Figure 3b: crosses); another source is from 107 flask
 361 data sites which observe CO₂ concentrations near the surface every week (Figure 3b: closed
 362 circles); lastly, GOSAT column data [Yokota *et al.*, 2004] are used (Figure 3b: gray lines),
 363 with orbital return periods of three days. For simplicity, in this simulation we did not
 364 account for the impact of cloud screening. We assume that the GOSAT data have the
 365 same averaging kernel as OCO [Wang *et al.*, 2009], i.e., nearly constant from the surface to
 366 the top of atmosphere. For this column data, the column observation increments are
 367 localized to each vertical level by the normalized averaging kernel for each level as
 368 follows:

$$369 \quad \mathbf{y}^b = h(\mathbf{x}^b) = \mathbf{A}^T (\mathbf{H}\mathbf{x}^b) = \sum_{i=1}^k a_i (\mathbf{H}\mathbf{x}_i^b) \quad (6)$$

370 where k is the number of vertical levels, \mathbf{H} the spatial interpolation operator, \mathbf{y}^b the model
 371 predicted CO₂ column mixing ratio, \mathbf{A} the averaging kernel, and a_i the element of \mathbf{A} at
 372 the i -th vertical level. We localize the j -th ensemble forecast column CO₂ to i -th vertical

373 level by the i -th averaging kernel element a_i as $y_{j,i}^b = a_i \times y_j^b$ and the column CO₂
374 observations to the i -th vertical level by a_i as $y_i^o = a_i \times y^o$. Then, $y_{j,i}^b$ and y_i^o are
375 compared during the analysis.

376 In the data assimilation system, the same model as the “nature run” is used for the
377 ensemble forecasts of 20 members (K=20), so that there is no model error (except for the
378 last experiment where the nature model has variable carbon fluxes not included in the
379 forecast model). Our goal is to test the impact of “variable localization” schemes in
380 estimating the spatial distribution of true CF shown in Figure 2a. Since CF is a forcing
381 term in the SPEEDY-C not changed by the forecast, it is updated only by the analysis step
382 of data assimilation, and the updated forcing from the analysis is then used for the next
383 forecast.

384 The initial ensemble members are chosen by random sampling from a long term
385 simulation of the SPEEDY-C and a SPEEDY-C coupled with a dynamic terrestrial carbon
386 model VEGAS [Zeng *et al.*, 2005] (hereafter referred as SPEEDY-VEGAS; Kang, 2009) in
387 order to generate fields of the initial ensemble background with no *a priori* information
388 about the nature run: 20 states of (U, V, T, q, Ps) and C are chosen randomly in time over an
389 one-year SPEEDY-C output and a three-year SPEEDY-VEGAS run respectively, and then

390 they are added by small random perturbations. For CF, from 20 fields of CO₂
391 concentration of the SPEEDY-C run in the midlevel at arbitrary times, we subtract the one-
392 day prior state of CO₂ concentration, and then convert the units of the field from ppmv/day
393 to kg/m²/s. Figures 2c and 2d show that the initial ensemble mean of the surface carbon
394 fluxes and the first level atmospheric CO₂ are very different from the true states in terms of
395 both spatial patterns and intensity. Since CO₂ concentration is well-mixed in the midlevel,
396 Figure 2d has very small values. Starting from these initial conditions without any *a*
397 *priori* information, we carried out the analyses of all dynamic variables for four months
398 using an analysis cycle of six hours.

399 The experimental settings described above are used for testing all schemes introduced
400 in this study to see the impact of “variable localization” techniques. In addition to these
401 experiments carried out in a perfect model and constant flux configuration, we have also
402 done another set of experiments testing the impact of variable carbon fluxes in the nature
403 model. With the same configuration of the observations and the same initial conditions,
404 we repeated the *L-Iway* and *C-univ* experiments now including terrestrial and oceanic CO₂
405 fluxes which evolve in time. We replace in “nature run” the CO₂ forcing every six hours
406 by the land surface CO₂ fluxes computed by VEGAS [Zeng *et al.* 2005] that includes the

407 vegetation impact on the carbon cycle, and the monthly prescribed oceanic fluxes
408 [Takahashi *et al.*, 2002] in addition to the fossil fuel emission used in the previous
409 experiments. We have produced one-year analysis and show the results for the last two-
410 month average in Section 4.

411

412 **4. Results**

413 Table 2 contains the global RMS errors for all variables from all the analysis schemes
414 during the last week of the 4-month data assimilation, and Figures 4 and 5 show the time
415 evolution of the global RMS errors in zonal wind and carbon variables. Other
416 meteorological variables have a similar pattern of RMS errors in the time series plot,
417 compared to the zonal wind. First, the standard fully multivariate data assimilation (*mult*)
418 has the worst results for all the variables. This is because *mult* allows for error
419 covariances among all variables in the analysis even though there is no physical
420 relationship between (C, CF) on the one hand and (T, q, Ps) on the other in the “nature”
421 model. Therefore, the estimations of the error covariances among these variables are only
422 due to sampling errors. Moreover, a poor representation of the initial surface carbon flux
423 can contaminate analyses of all variables in *mult*. As a result, the *mult* analysis has larger

424 errors and eventually undergoes filter divergence, i.e., the feedback from the sampling
425 errors makes the analysis of meteorological variables so poor that the diagnosis of the
426 model variables in the forecast fails after 40 days. In theory, this problem of the *mult*
427 system could be resolved by using much larger ensemble size so that sampling errors are
428 reduced, but in practice this approach is not computationally feasible.

429 By eliminating the unphysical relationship between the carbon flux CF and (U, V, T, q,
430 Ps), *L-mult* prevents a poor initial representation of CF from degrading the analysis of the
431 other variables. Also, the analysis of carbon variables benefits from better states of other
432 variables (without a contamination of the surface carbon flux). As a result, the *L-mult*
433 analysis is improved significantly for all dynamic variables and filter divergence is avoided.
434 Still, there is unnecessary feedback between C and (T, q, Ps), which is negligible in nature.
435 Thus, *L-mult* is not the optimal method and can be improved further by additional variable
436 localization (Figure 4, 5 and Table 2).

437 In *Iway*, we zero out the background error covariance between (C, CF) and (T, q, Ps).
438 Furthermore, *Iway* does not allow any changes in the meteorological variables due to the
439 CO₂ variables. Compared with *mult*, this does not allow any feedback between carbon
440 variables (C, CF) and (T, q, Ps), but, in contrast with *L-mult*, it does include the covariance

441 between CF and wind fields. Carbon variables from *Iway* analysis are improved
442 significantly while the analyses of meteorological variables are, as expected, comparable
443 with the results of *L-mult*. Mean RMS errors (Table 2) show that the differences between
444 *Iway* and *L-mult* are only on the order of 1% for (U, V, T, q, Ps) whereas *Iway* improves
445 the estimates of (C, CF) by 30-35%.

446 Figure 6 compares maps of the analysis errors in the zonal and meridional wind fields
447 obtained with *Iway* and with *L-mult*. Due to the distribution of the rawinsonde network
448 sites (Figure 3a), errors are large over the oceans and polar regions in both *L-mult* and
449 *Iway*. The analysis of wind in *Iway* has similar error patterns but smaller error
450 amplitudes than in *L-mult*. By contrast, *Iway* results in a major improvement in the
451 atmospheric CO₂ analysis as shown in Figure 7. Since *L-mult* considers the background
452 of (T, q, Ps) in addition to (U, V, C, CF) for analyzing the atmospheric CO₂, the
453 background uncertainties of (T, q, Ps) can influence the weight between the background
454 and the observations of atmospheric CO₂. Although a large uncertainty of temperature
455 can be related to the wind uncertainty so that the carbon dioxide concentration could be
456 affected by those wind errors, this is not a first-order effect and does not need to be
457 considered during the analysis. This is what the result in Figure 7 shows: *L-mult* has

458 larger errors overall and the spatial pattern is not as smooth as the nature run or the results
459 from *Iway*. Because the *L-mult* analysis reflects more strongly the observations of
460 atmospheric carbon whenever there are large background uncertainties of (T, q, Ps) in
461 addition to (U, V, C, CF), atmospheric CO₂ observations are over-weighted for the case of
462 *L-mult* producing an analysis with additional noise (Figure 7a) compared with the case of
463 *Iway* (Figure 7c). Over the ocean, where there are few observations of meteorological
464 variables, their estimated error, given by the background spread, is large. Thus, the *L-*
465 *mult* tends to give more weights to the atmospheric CO₂ observations than it should
466 because it considers the joint background uncertainties of (U, V, T, q, C, Ps) altogether.

467 We further localize the variables in *L-Iway* by zeroing out the correlation between CF
468 and (U,V) from the *Iway* system. The analysis can still include the uncertainties in the
469 wind field to assist the analysis of atmospheric CO₂, but the error of surface carbon flux is
470 coupled with only the atmospheric CO₂ uncertainty reflecting the fact that carbon flux is
471 only related to low level atmospheric CO₂ and not with the wind. Again, the
472 meteorological variables are not affected by (C, CF), so that the analysis of (U, V, T, q, Ps)
473 are exactly the same as in *Iway*, also true for the *C-univ* analysis for the same reason.
474 From Table 2 and Figure 5, we find that *L-Iway* has the best performance of five schemes

475 for estimating surface CO₂ fluxes, while the result for atmospheric CO₂ concentration is
476 comparable with that from *Iway* (Figure 4). This implies that the surface carbon fluxes
477 should not be linked to the wind fields in the background error covariance matrix. As a
478 result, the spatial distribution of the analysis from *L-Iway* in Figure 8 also shows a
479 promising performance in estimating surface carbon fluxes, capturing well the major source
480 regions in the Northern Hemisphere.

481 The last method considered, *C-univ*, has stable results in the analysis of the carbon
482 variables, but the surface carbon flux is slightly worse than that of *L-Iway* (Figure 4, 5 and
483 Table 2). Interestingly, the RMS error of surface carbon analysis grows with time whereas
484 *L-Iway* keeps reducing the errors (Figure 5). Since these two systems differ only in
485 whether the transport error is considered when analyzing the atmospheric CO₂
486 concentrations, the gradual increase of RMS error in *C-univ* can be seen as a result of
487 neglecting transport errors.

488 Figure 9 displays global maps of analysis errors in surface CO₂ flux analyses resulting
489 from the *L-mult*, *Iway*, *L-Iway*, and *C-univ* experiments (recall that the standard
490 multivariate LETKF without any variable localization blew up after 40 days). *L-mult* has
491 a broad area of overall errors (Fig 9a). It is apparent that the presence of an error

492 covariance among all of the atmospheric variables is not helpful for the analysis of carbon,
493 since it just introduces sampling errors. By removing the irrelevant error covariance
494 between carbon and temperature, humidity, and surface pressure from *L-mult*, the results in
495 *Iway* show improvement overall (Fig 9b) compared to the multivariate analyses. *L-Iway*
496 provides further localization between the surface carbon flux and wind fields, compared to
497 *Iway*, and hence obtains the smallest errors in carbon flux analysis. This technique clearly
498 has less error, especially over the oceans, than *L-mult* or *Iway*.

499 The approach embodied in *C-univ* has lost the error information contained in the
500 relationship between wind and atmospheric CO₂ uncertainties and hence has somewhat
501 worse results than *L-Iway*. Indeed, over the polar regions (Figure 9d), *C-univ* has
502 spurious estimates of surface carbon fluxes in areas where there are large errors in the wind
503 analysis (Figure 6), whereas *L-Iway* does not have those errors. We also note that the
504 error in *C-univ* over the polar region grows with time, and this leads to RMS error
505 increases in Figure 5. Thus, we can conclude that the reason for increasing RMS error in
506 the surface carbon fluxes is that transport errors are not accounted for in *C-univ*. In
507 addition, experiments with an imperfect model [Kang, 2009] indicate that the perfect model
508 assumption underestimates the impact of this deficiency of *C-univ* since transport errors are

509 also underestimated in this scenario.

510 When we allow for time-varying surface CO₂ forcing, the estimation problem becomes
511 more difficult because we are not anymore under a perfect model scenario, since the
512 forecast model does not change the surface fluxes, only changed by the analysis cycle.
513 Thus, the overall errors for both schemes become larger and require further research on
514 potential improvements in the data assimilation techniques (see below). Nevertheless, the
515 relatively small advantage of *L-Iway* compared to *C-univ* observed with constant fluxes
516 (Fig. 9c and Fig 9d) becomes much more significant (Figure 10) indicating that, for this
517 scenario, the estimation of surface CO₂ fluxes from *L-Iway* is significantly better than that
518 from *C-univ*. It is important to note that *L-Iway* outperforms *C-univ* especially over the
519 ocean and the Southern Hemisphere where the wind uncertainties are dominant due to the
520 lack of rawinsonde observations. Since the analysis cycle updates surface CO₂ fluxes
521 which in turn force the atmospheric CO₂ forecast for the next six hours, unresolved
522 transport errors when assimilating CO₂ in *C-univ* can degrade the analysis of carbon
523 variables more in the case with the time-varying forcing than in the case with a constant
524 forcing.

525 We note that the adaptive inflation estimation has relatively large changes during the

526 first ten days of the analysis when the errors in the initial conditions of the background
527 states are very large compared to the observation errors (not shown). The adaptive
528 inflation of the background covariance for the meteorological variables, which is estimated
529 initially to be about 35%, settles after spin-up at about 5% ($\rho \approx 1.05$ in Equation 4). The
530 inflation factor estimated for the atmospheric CO₂ concentration also decreases with time:
531 the inflation factor is estimated at about 50% during the first week and then converges in
532 time to less than 10%. These adaptive inflation factors are similar for all the variable
533 localization schemes that we have examined in this study. The inflation for the surface
534 carbon fluxes is estimated to be small, less than 2%, as could be expected for a variable that
535 is not observed [Anderson, 2009]. If instead, we allow the inflation for the carbon flux to
536 be the same as for atmospheric CO₂, there is filter divergence in the estimation of the
537 surface carbon flux analysis. Thus, the adaptive inflation estimation algorithm [Li *et al.*,
538 2009; Zhang *et al.*, 2004] appears to work quite well in the carbon cycle data assimilation
539 system.

540

541 **5. Summary and Discussion**

542 We have developed a method to estimate surface carbon fluxes via an EnKF data

543 assimilation analyzing the meteorological variables and the carbon variables
544 simultaneously. The method is fairly efficient in terms of computational cost since it does
545 not require an additional run of the transport model as the observation operator during the
546 analysis, a step generally used in previous studies. In addition, simultaneous analyses
547 allow accounting for the important day-to-day wind uncertainties when analyzing CO₂
548 variables. Atmospheric CO₂ observations are assimilated from a simulated network of *in*
549 *situ* (continuous record), flask (weekly record), and satellite-based measurements with
550 realistic resolution. The results of this study, although far from perfect, are promising
551 especially considering that no *a priori* information about carbon has been used.

552 The focus of this paper is a comparison of several “variable localization” schemes that
553 reduce sampling errors in the ensemble estimation of the covariance between physically
554 uncorrelated variables by zeroing out the background error covariance among these
555 variables. Since carbon variables in the nature run do not have a physical relationship
556 with temperature, specific humidity and surface pressure, the standard EnKF approach of
557 coupling errors of all variables in *mult* analysis induces sampling error into the system.
558 As a result, the accuracy of *mult* analyses for all dynamic variables gets progressively
559 worse together with surface carbon flux estimation until about 40 days, when filter

560 divergence takes place. Of the five new methods introduced here, the localized one-way
561 approach, *L-Iway*, has the best performance in the estimation of surface carbon fluxes. The
562 atmospheric CO₂ analysis includes the error covariance of CO₂ and surface carbon flux as
563 well as the wind transport error, which is strongly related to the forecast of atmospheric
564 CO₂. This approach excludes the non-physical error covariance between the wind field
565 and surface CO₂ flux and among the carbon variables and temperature, humidity, and
566 surface pressure, which are dominated by sampling errors. Moreover, the carbon variables
567 are not allowed to influence the analysis of meteorological variables because CO₂ is poorly
568 observed, and thus would increase the sampling errors in the better observed winds and
569 temperatures [Liu et al., 2009].

570 The results from *L-Iway* can be contrasted with *C-univ*, which is closer to previous
571 studies in a sense that transport error covariances are not considered during the carbon
572 analysis. Nevertheless, the *C-univ* approach within EnKF does allow for information on
573 transport uncertainties because the different ensemble members have different winds, and
574 therefore different CO₂ transports. As a result, the carbon univariate approach gives quite
575 good results when we use constant surface fluxes, although slightly worse than those
576 obtained with the *L-Iway* approach. The improvement of *L-Iway* over *C-univ* becomes

577 much larger when the imperfection of CO₂ forecast becomes important. The advantages
578 of *L-Iway* results compared to *C-univ* results demonstrate that it is necessary to resolve
579 transport error for the analysis of atmospheric CO₂.

580 We note that the variable localization design of the most successful method in this
581 paper, the *L-Iway*, is based on our OSSE experimental setting since, in our nature run,
582 atmospheric CO₂ is only transported and mixed by the wind fields and the varying CO₂ has
583 no radiative impact and thus no temperature dependence. In a more realistic model,
584 assimilating real observations, the variable localization technique we have introduced needs
585 to be adapted by considering the “*dynamical distance*” between each pair of variables in a
586 real nature and model. If the background error covariance is dominated by sampling
587 errors, it will be beneficial to zero out the covariance as we did here, even if the two
588 variables are, to some extent, physically related. For example, biospheric and air-sea
589 carbon fluxes have diurnal, seasonal, and interannual variabilities that are modulated by
590 precipitation, temperature, cloud cover, relative humidity, and wind speeds. Only if the
591 atmospheric carbon model is realistic enough to represent well the covariability of two of
592 these variables, should the corresponding error covariance be retained. Furthermore, a
593 study with more realistic settings such as using a realistic model and an imperfect model

594 assumption is required as a next step, in order to further examine the impact of this new
595 method on assimilating real observations.

596 We point out that, in this paper, we introduced the methodology of constraining the
597 unobserved surface CO₂ fluxes by assimilating atmospheric CO₂ observations
598 simultaneously with atmospheric observations allowing transport errors to be considered
599 during the analysis step. In principle, this methodology could be extended to the
600 estimation of surface moisture/heat fluxes from the assimilation of observations of
601 humidity/temperature in the atmosphere, another major challenge in current models.

602 Finally, we note that the results of these new techniques such as variable localization
603 and adaptive inflation have clearly improved our ability to estimate the surface fluxes, so
604 that these techniques can be used in other Ensemble Kalman Filter data assimilation
605 problems. Nevertheless, since our ultimate goal is to estimate as well as possible not only
606 the atmospheric CO₂ but also the surface carbon fluxes, it is clear that significant more
607 progress is needed, especially in the imperfect model scenario. We are doing research
608 with several promising additional new techniques, including the estimation of the model
609 bias, and the restructuring of ensemble perturbations that in time tend to align themselves
610 too much along the most unstable direction (leading local Lyapunov vectors). The

611 difficulty of the problem makes clear the need to perform OSSEs as well as real data
612 experiments in order to understand what can be achieved with real data and what
613 techniques should be tested.

614 **Acknowledgments**

615 We are grateful to the US Department of Energy for the support of the research project,
616 “Carbon data assimilation with coupled Ensemble Kalman filter”, under DOE Grant
617 DEFG0207ER64437. Support was also received from NASA Grants NNX08AD4oG,
618 NNX07AM97G, NOAA Grant NA09OAR4310178, and ONR Grant N000141010557.
619 The SPEEDY model was kindly provided by Franco Molteni and Fred Kucharski. The
620 very constructive suggestions of Andy Jacobson and two anonymous reviewers improved
621 the paper content and presentation.

622 **References**

623

624 Anderson J. L. (2001), An Ensemble Adjustment Kalman Filter for Data Assimilation. *Mon.*
625 *Wea. Rev.* *129*, 2884-2903.

626 Anderson, J. L. (2009), Spatially and temporally varying adaptive covariance inflation for
627 ensemble filters. *Tellus*, *61A*, 72–83.

628 Anderson, J. L. and S.L. Anderson (1999), A Monte Carlo implementation of the nonlinear
629 filtering problem to produce ensemble assimilations and forecasts. *Mon. Wea. Rev.* *127*,
630 2741–2758.

631 Andres, R. J., G. Marland, I. Fung, and E. Matthews (1996), A 1° x 1° distribution of carbon
632 dioxide emissions from fossil fuel consumption and cement manufacture, 1950-1990,
633 *Global Biogeochem. Cycles*, *10*, 419-429.

634 Annan, J. D., D. J. Lunt, J. C. Hargreaves, and P. J. Valdes (2005), Parameter estimation in
635 an atmospheric GCM. *Nonlinear processes in geophysics*, **12**, 363–371.

636 Baek, S.-J., B.R. Hunt, E. Kalnay, E. Ott, I. Szunyogh (2006), Local ensemble Kalman
637 filtering in the presence of model bias, *Tellus*, *58A*, 293-306.

638 Baker, D. F., S. C. Doney, and D. S. Schimel (2006), Variational data assimilation for

639 atmospheric CO₂, *Tellus*, 58B, 359-365

640 Baker, D. F., Bösch, H., Doney, S. C., O'Brien, D., and Schimel, D. S.: Carbon source/sink
641 information provided by column CO₂ measurements from the Orbiting Carbon
642 Observatory, *Atmos. Chem. Phys.*, 10, 4145-4165, doi:10.5194/acp-10-4145-2010.

643 Bishop, C. H., B. J. Etherton, and S. J. Majumdar (2001), Adaptive Sampling with the
644 Ensemble Transform Kalman Filter. Part I: Theoretical Aspects. *Mon. Wea. Rev.* 129,
645 420-436.

646 Bousquet, P., P. Peylin, P. Ciais, C. Le Quere, P. Friedlingstein, and P. P. Tans (2000),
647 Regional changes in carbon dioxide fluxes of land and oceans since 1980. *Science*,
648 290, 1342-1346.

649 Chevallier, F., R. J. Engelen, C. Carouge, T. J. Conway, P. Peylin, C. Pickett-Heaps, M.
650 Ramonet, P. J. Rayner, and I. Xueref-Remy (2009a), AIRS-based versus flask-based
651 estimation of carbon surface fluxes, *J. Geophys. Res.*, 114, D20303,
652 doi:10.1029/2009JD012311.

653 Chevallier, F., S. Maksyutov, P. Bousquet, F.-M. Bréon, R. Saito, Y. Yoshida, and T. Yokota
654 (2009b), On the accuracy of the CO₂ surface fluxes to be estimated from the GOSAT
655 observations, *Geophys. Res. Lett.*, 36, L19807, doi:10.1029/2009GL040108.

656 Crisp, D., R. M. Atlas, F.-M. Breon, L. R. Brown, J. P. Burrows, P. Ciais, B. J. Connor, S.
657 C. Doney, I. Y. Fung, D. J. Jacob, E. Miller D. O'Brien, S. Pawson, J. T. Randerson,
658 P. Rayner, R. J. Salawitch, S. P. Sander, B. Sen, G. L. Stephens, P. P. Tans, G. C.
659 Toon, P. O. Wennberg, S. C. Wofsy, Y. L. Yung, Z. Kuang, B. Chudasama, G.
660 Sprague, B. Weiss, R. Pollock, D. Kenyon, and S. Schroll (2004), The Orbiting
661 Carbon Observatory (OCO) mission. *Advances in Space Research*, 34, 700-709
662 Desroziers G., L. Berre, B. Chapnik, and P. Poli (2005), Diagnosis of observation,
663 background and analysis error statistics in observation space. *Quart. J. Roy. Meteor.*
664 *Soc.*, 131, 3385-3396.
665 Engelen, R. J., S. Serrar, and F. Chevallier (2009), Four-dimensional data assimilation of
666 atmospheric CO₂ using AIRS observations, *J. Geophys. Res.*, 114, D03303,
667 doi:10.1029/2008JD010739.
668 Enting, I. G. (2002), *Inverse Problems in Atmospheric Constituent Transport*, Cambridge
669 University Press, N. Y.
670 Evensen, G. (1994), Sequential data assimilation with a nonlinear quasi-geostrophic model
671 using Monte Carlo methods to forecast error statistics. *J. Geophys. Res.*, 99, 10 143-
672 10 162.

673 Feng, L., Palmer, P. I., Bösch, H., and Dance, S.: Estimating surface CO₂ fluxes from
674 space-borne CO₂ dry air mole fraction observations using an ensemble Kalman Filter,
675 *Atmos. Chem. Phys.*, **9**, 2619-2633, doi:10.5194/acp-9-2619-2009, 2009.

676 Gaspari, G., and S. E. Cohn (1999), Construction of correlation functions in two and three
677 dimensions. *Quart. J. Roy. Meteor. Soc.*, **125**, 723-757.

678 Greybush, S., E. Kalnay, T. Miyoshi, K. Ide, and B. Hunt (2010), Balance and Ensemble
679 Kalman Filter Localization Techniques, *Mon. Wea. Rev.*, in press, doi:
680 10.1175/2010MWR3328.1.

681 Gurney, K. R., R. M. Law, A. S. Denning, P. J. Rayner, B. C. Pak, D. Baker, P. Bousquet, L.
682 Bruhwiler, Y. H. Chen, P. Ciais, I. Y. Fung, M. Heimann, J. John, T. Maki, S.
683 Maksyutov, P. Peylin, M. Prather, and S. Taguchi (2004), Transcom 3 inversion
684 intercomparison: Model mean results for the estimation of seasonal carbon sources
685 and sinks. *Global Biogeochemical Cycles*, **18**(1), 10.1029/2003GB002111.

686 Hamill, T. M., J. S. Whitaker, and C. Snyder (2001), Distance-Dependent Filtering of
687 Background Error Covariance Estimates in an Ensemble Kalman Filter. *Mon. Wea.*
688 *Rev.*, **129**, 2776-2789.

689 Hollingsworth, A., et al. (2008), The Global Earth-system Monitoring using Satellite and

690 in-situ data (GEMS) Project: Towards a monitoring and forecasting system for
691 atmospheric composition, *Bull. Am. Meteorol. Soc.*, 89, 1147–1164,
692 doi:10.1175/2008BAMS2355.1.

693 Houtekamer, P. L. and H. L. Mitchell (2001), A Sequential Ensemble Kalman Filter for
694 Atmospheric Data Assimilation. *Monthly Weather Review*: Vol. 129, pp. 123–137.

695 Hunt, B. R., E. Kostelich, and I. Szunyogh (2007), Efficient Data Assimilation for
696 Spatiotemporal Chaos: a Local Ensemble Transform Kalman Filter, *Physica D*, 230,
697 112-126.

698 Kang, J. (2009), Carbon Cycle Data Assimilation Using a Coupled Atmosphere-Vegetation
699 Model and the Local Ensemble Transform Kalman Filter, Ph.D Thesis, University of
700 Maryland

701 Li, H., E. Kalnay, and T. Miyoshi, (2009), Simultaneous estimation of covariance inflation
702 and observation errors within ensemble Kalman filter. *Quart. J. Roy. Meteor. Soc.*,
703 135, 523-533.

704 Liu, J., E. Kalnay, I. Fung, M. Chahine, and E. T. Olsen (2011), Assimilation of AIRS CO2
705 observations with an EnKF in a Carbon-Climate Model, paper presented at 91st
706 American Meteorological Society Annual Meeting, Seattle, WA, USA.

707 Liu, J., H. Li, E. Kalnay, E.J. Kostelich, and I. Szunyogh (2009), Univariate and
708 Multivariate Assimilation of AIRS Humidity Retrievals with the Local Ensemble
709 Transform Kalman Filter. *Mon. Wea. Rev.*, *137*, 3918–3932.

710 Miyazaki, K. (2009), Performance of a local ensemble transform Kalman filter for the
711 analysis of atmospheric circulation and distribution of long-lived tracers under
712 idealized conditions, *J. Geophys. Res.*, *114*, D19304, doi:10.1029/2009JD011892.

713 Miyoshi, T (2005), Ensemble Kallman Filter Experiments with a Primitive-Equation Global
714 Model, Ph.D Thesis, University of Maryland

715 Molteni, F. (2003), Atmospheric simulations using a GCM with simplified physical
716 parametrizations. I: Model climatology and variability in multi-decadal experiments.
717 *Climate Dyn.*, *20*, 175-191.

718 Ott, E., B. R. Hunt, I. Szunyogh, A. V. Zimin, E. J. Kostelich, M., Corazza, E. Kalnay, D. J.
719 Patil, and J. A. Yorke (2004), Estimating the state of large spatio- temporally chaotic
720 systems. *Phys. Lett. A.*, *330*, 365- 370.

721 Peters, W., J. B. Miller, J. Whitaker, A. S. Denning, A. Hirsch, M. C. Krol, D. Zupanski, L.
722 Bruhwiler, and P.P. Tans (2005), An ensemble data assimilation system to estimate
723 CO₂ surface fluxes from atmospheric trace gas observations, *J. Geophys. Res.*, *110*,

724 D24304, doi:10.1029/2005JD006157.

725 Peters, W., et al. (2007), An atmospheric perspective on North American carbon dioxide
726 exchange: Carbon Tracker, *Proc. Natl. Acad. Sci. USA.*, *104*, 18,925– 18,930.

727 Rödenbeck, C., S. Houweling, M. Gloor, M. Heimann (2003), CO₂ flux history 1982-2001
728 inferred from atmospheric data using a global inversion of atmospheric transport,
729 *Atmos. Chem. Phys.*, *3*, 1919-1964.

730 Stephens, B. B., et al. (2007), Weak northern and strong tropical land carbon ptake from
731 vertical profiles of atmospheric CO₂, *Science*, *316*, 1732– 1735.

732 Takahashi, T., S. C. Sutherland, C. Sweeney, A. Poisson, N. Metzl, B. Tilbrook, N. Bates,
733 R. Wanninkhof, R. A. Feely, C. Sabine, J. Olafsson, Y. Nojiri, 2002: Global sea-air
734 CO₂ flux based on climatological surface ocean pCO₂, and seasonal biological and
735 temperature effects, *Deep-Sea Research II*, *49*, 1601-1622.

736 Wang, H., D. J. Jacob, M. Kopacz, D. B. A. Jones, P. Suntharalingam, J. A. Fisher, R.
737 Nassar, S. Pawson, and J. E. Nielsen (2009), Error correlation between CO₂ and
738 CO as constraint for CO₂ flux inversions using satellite data, *Atmos. Chem. Phys.*, *9*,
739 7313-7323.

740 Yokota, T., H. Oguma, I. Morino, and G. Inoue (2004), A nadir looking SWIR FTS to

741 monitor CO₂ column density for Japanese GOSAT project, *Proc. Twenty-fourth Int.*
742 *Sympo. on Space Technol. and Sci (Selected Papers), JSASS and Organizing Comm.*
743 *of the 24th ISTS*, 887–889.

744 Zeng, N., A. Mariotti, and P. Wetzel, 2005: Terrestrial mechanisms of interannual CO₂
745 variability, *Global Biogeochemical Cycles*, 19, GB1016, doi:10.1029/2004GB002273.

746 Zhang F., C. Snyder, and J. Sun (2004), Impacts of initial estimate and observation
747 availability on convective-scale data assimilation with an ensemble Kalman filter.
748 *Mon. Wea. Rev.*, 132, 1238–1253.

749 Zupanski, D., A. S. Denning, M. Uliasz, M. Zupanski, A. E. Schuh, P. J. Rayner, W. Peters,
750 and K. D. Corbin (2007), Carbon flux bias estimation employing Maximum
751 Likelihood Ensemble Filter (MLEF), *J. Geophys. Res.*, 112, D17107,
752 doi:10.1029/2006JD008371.

753 Zupanski, M., 2005: Maximum Likelihood Ensemble Filter: Theoretical Aspects. *Mon.*
754 *Wea. Rev.*, 133, 1710–1726. doi: 10.1175/MWR2946.1

755 Zupanski, M., S. J. Fletcher, I. M. Navon, B. Uzunoglu, R. P. Heikes, D. A. Randall, T. D.
756 Ringlee and D. Daescu, 2006: Initiation of ensemble data assimilation. *Tellus*, 58A,
757 159-170.

758 **Figure Captions**

759 Figure 1. Schematic plots of background error covariance matrices ($P^b=(x^b)(x^b)^T/(K-1)$) for
760 (a) *mult*, (b) *L-mult*, (c) *1way*, (d) *L-1way*, and (e) *C-univ* analysis systems. Here, C
761 indicates atmospheric CO₂ concentration and CF indicates surface carbon fluxes. The
762 colors of the variable names are matched with the system used for their updates. White
763 areas with “no” indicate the error correlation between variables is assumed to be zero
764 during the analysis while areas with “yes” indicate that the errors are allowed to be
765 correlated. For example, in 1(d), the errors of the standard atmospheric variables are
766 coupled, the atmospheric CO₂ errors are coupled with those of the wind but the wind errors
767 are not coupled with the CO₂ errors (1-way coupling), and the surface carbon flux errors are
768 only coupled with the CO₂ errors.

769

770 Figure 2. True state of (a) surface CO₂ fluxes (6 PgC/yr) and (b) atmospheric CO₂
771 concentrations in the lowest layer at the initial time, as well as initial ensemble mean of (c)
772 surface CO₂ fluxes and (d) atmospheric CO₂. Units for atmospheric CO₂ concentration are
773 ppmv, units for surface CO₂ fluxes are 10⁻⁹ kg/m²/s.

774 Figure 3. The simulated observational coverage of (a) meteorological variables (black dots)
775 and (b) atmospheric CO₂ concentration (gray lines: GOSAT column data, crosses:
776 continuous *in situ* data, closed circles: weekly flask data).

777

778 Figure 4. Time series of global RMS error of (a) U (m/s) and (b) atmospheric CO₂
779 concentration in the lowest layer (ppmv) for four months of analysis. (solid gray: *mult*,
780 solid black: *L-mult*, dashed gray: *Iway*, dashed black: *L-Iway*, dotted light gray: *C-univ*)

781

782 Figure 5. Same as Figure 4, except for the surface CO₂ fluxes.

783

784 Figure 6. Analysis error (unit: m/s) of (a) zonal wind and (b) meridional wind from the
785 localized multivariate analysis (*L-mult*) for the last three months of data assimilation. (c)
786 and (d): The same as in (a) and (b) except from *Iway*. Shading indicates positive errors and
787 contours indicate negative errors with the same color scale as the shading.

788

789 Figure 7. Analysis (left column) of atmospheric CO₂ concentration in the lowest layer and
790 its error (right column) after four months of analysis. (a) and (b) results from *L-mult*, (c)
791 and (d) from *Iway*. Units are ppmv. Shading indicates positive errors and contours
792 indicate negative errors with the same color scale as the shading.

793

794 Figure 8. (a) True state of surface CO₂ fluxes and (b) the analysis after four months of the
795 L-1way (localized 1-way multivariate) data assimilation. (Units are 10⁻⁹ kg/m²/s.)

796

797 Figure 9. Analysis errors of surface CO₂ fluxes after four months of analysis. (a) results
798 from *L-mult*, (b) from *Iway*, (c) from *L-Iway* and (d) from *C-univ*. Units are 10⁻⁹ kg/m²/s.
799 (Shading indicates positive errors and contours indicate negative errors with the same color
800 scale as the shading.)

801

802 Figure 10. (a) True state of surface CO₂ fluxes from a time-varying terrestrial and oceanic
803 forcing and a fossil fuel emission, and the estimated surface CO₂ fluxes from (b) *L-Iway*,
804 and (c) *C-univ* data assimilation for the last two months (November-December) of one-year
805 analysis

806 Table 1. Standard deviation of errors used in creating the simulated observations.

807

Variable	Std. dev. of error
U	1.0 m/s
V	1.0 m/s
T	1.0 K
q	0.1 g/kg
Ps	1.0 hPa
C	1.0 ppmv

808 Table 2. RMS error of variables from different localization schemes for the last week of
 809 four-month analysis: one-week time average of every six hour values of

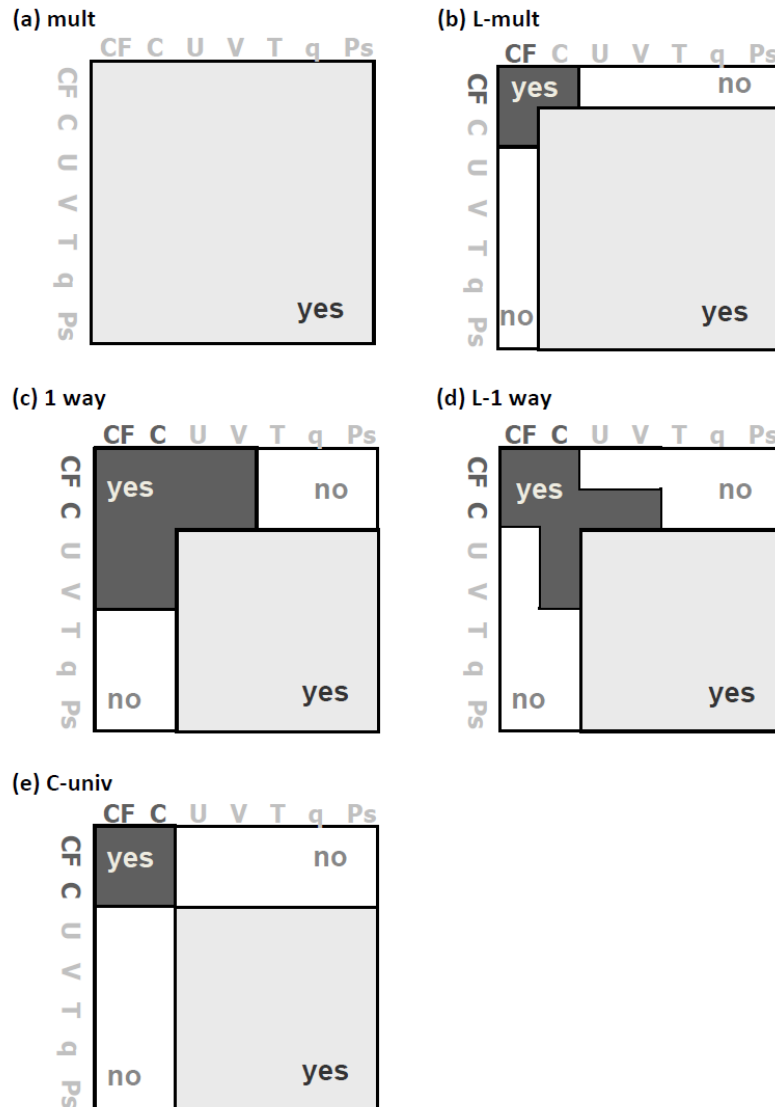
810 $\sqrt{\sum_i^n (x_i^a - x_i^t)^2} / n$, where x_i^a / x_i^t is the analysis/the truth at i-th point, and n is the total

811 number of grid points (units: CF= 10^{-9} kg/m²/s, C=ppmv, U and V=m/s, T=K, q=g/kg,
 812 Ps=hPa). The errors corresponding to Ensemble Kalman Filter divergence are symbolically
 813 represented as “infinite”.

814

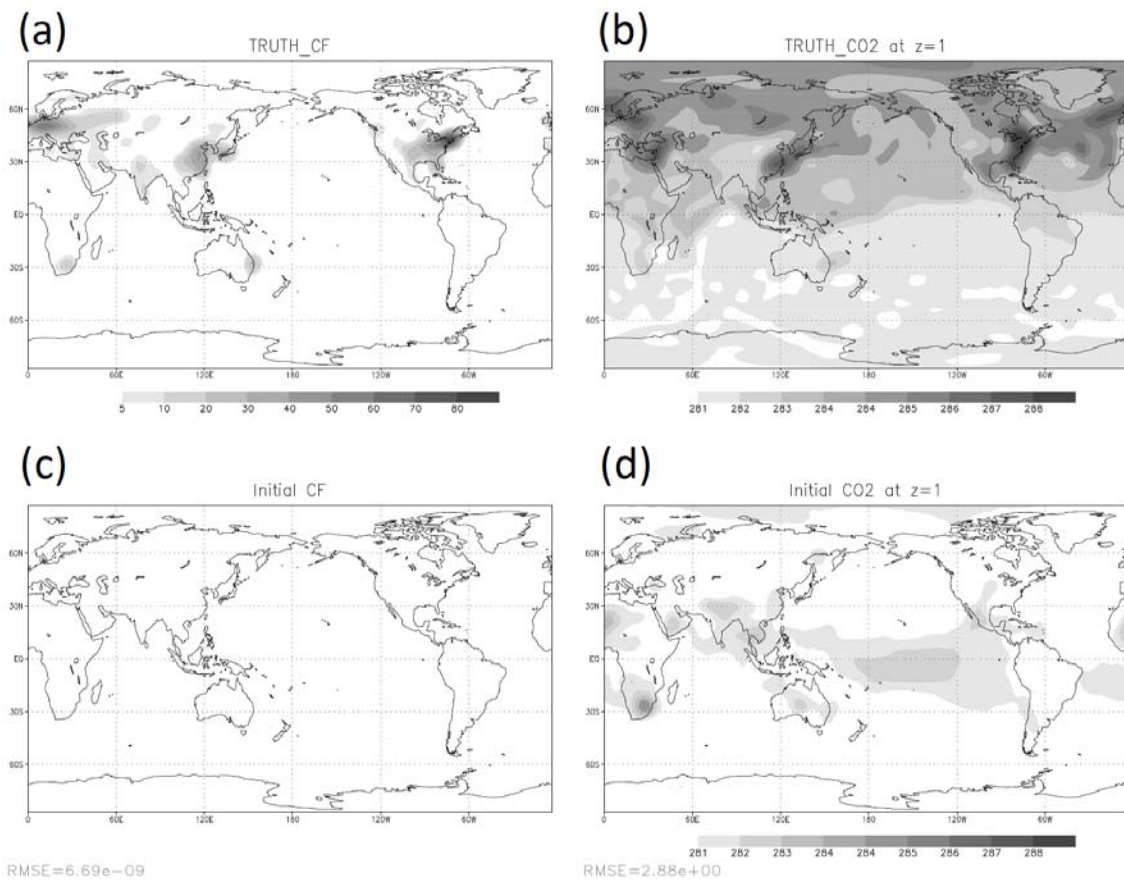
	<i>mult</i>	<i>L-mult</i>	<i>lway</i>	<i>L-lway</i>	<i>C-univ</i>
CF	∞	8.65	6.10	5.65	5.79
C	∞	1.05	0.68	0.71	0.67
U	∞	1.32		1.32	
V	∞	1.22		1.20	
T	∞	0.53		0.54	
q	∞	0.34		0.35	
Ps	∞	1.15		1.14	

815



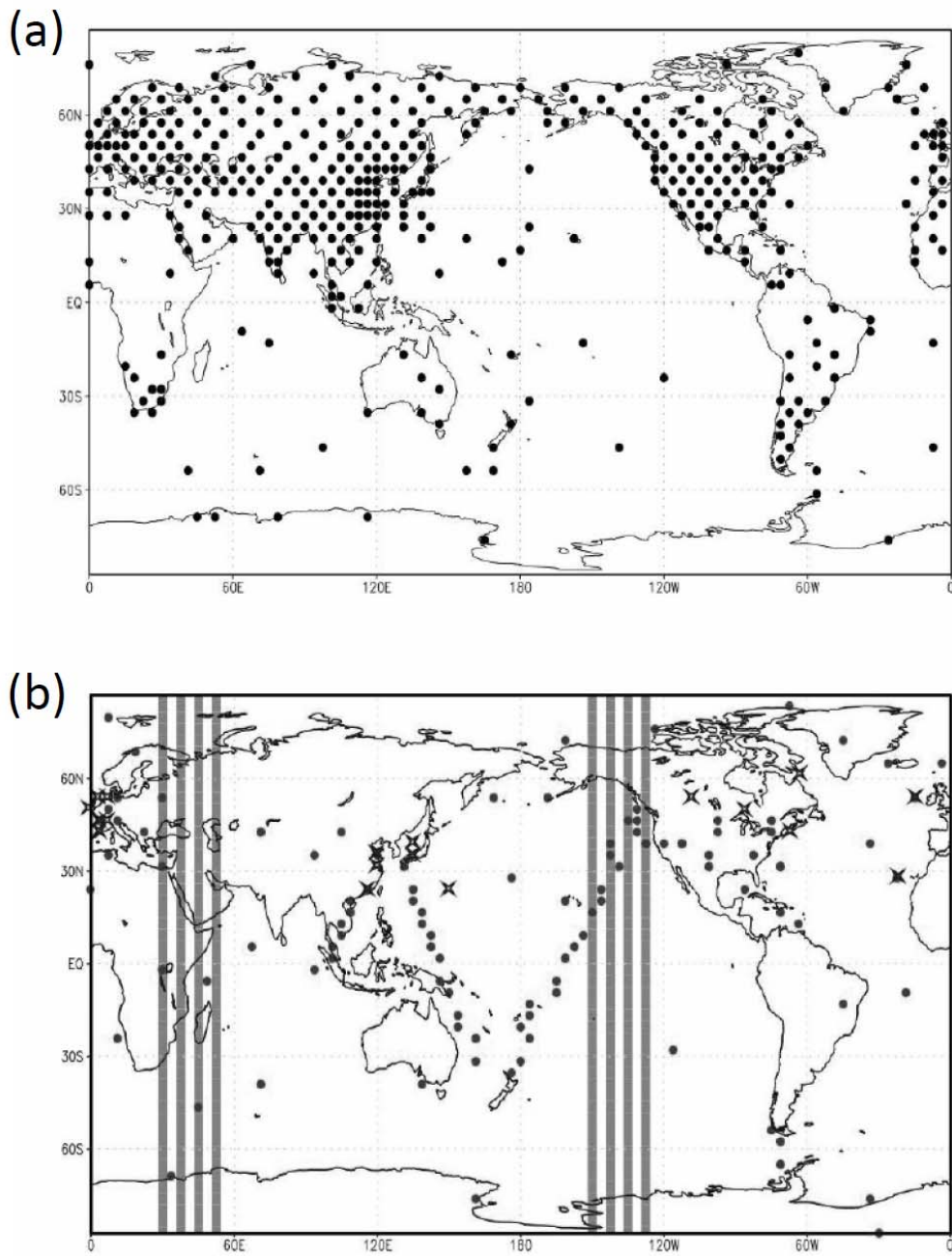
816

817 Figure 1. Schematic plots of background error covariance matrices ($P^b=(x^b)(x^b)^T/(K-1)$) for (a) *mult*, (b) *L-*
818 *mult*, (c) *1way*, (d) *L-1way*, and (e) *C-univ* analysis systems. Here, C indicates atmospheric CO₂
819 concentration and CF indicates surface carbon fluxes. The colors of the variable names are matched with the
820 system used for their updates. White areas with “no” indicate the error correlation between variables is
821 assumed to be zero during the analysis while areas with “yes” indicate that the errors are allowed to be
822 correlated. For example, in 1(d), the errors of the standard atmospheric variables are coupled, the atmospheric
823 CO₂ errors are coupled with those of the wind but the wind errors are not coupled with the CO₂ errors (1-way
824 coupling), and the surface carbon flux errors are only coupled with the CO₂ errors.



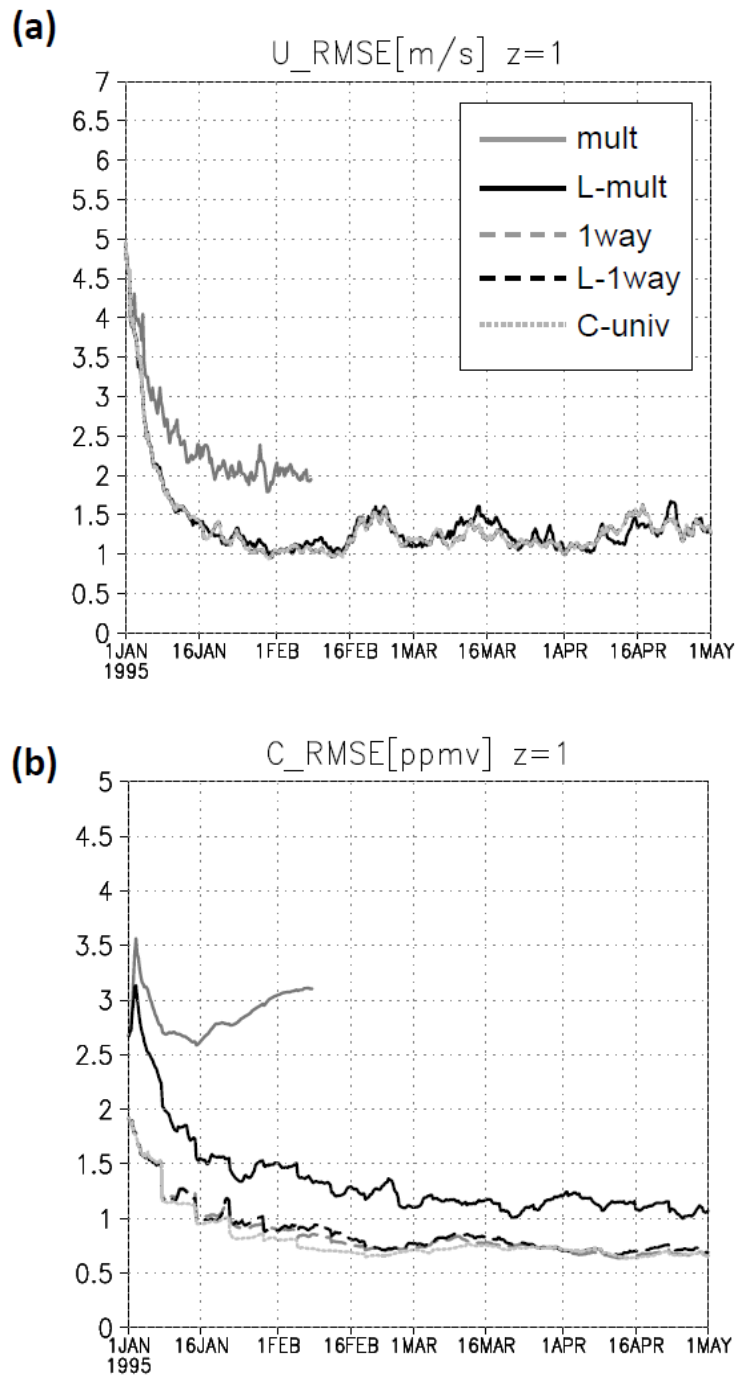
825

826 Figure 2. True state of (a) surface CO₂ fluxes (6 PgC/yr) and (b) atmospheric CO₂
 827 concentrations in the lowest layer at the initial time, as well as initial ensemble mean of (c)
 828 surface CO₂ fluxes and (d) atmospheric CO₂. Units for atmospheric CO₂ concentration are
 829 ppmv, units for surface CO₂ fluxes are 10⁻⁹ kg/m²/s.



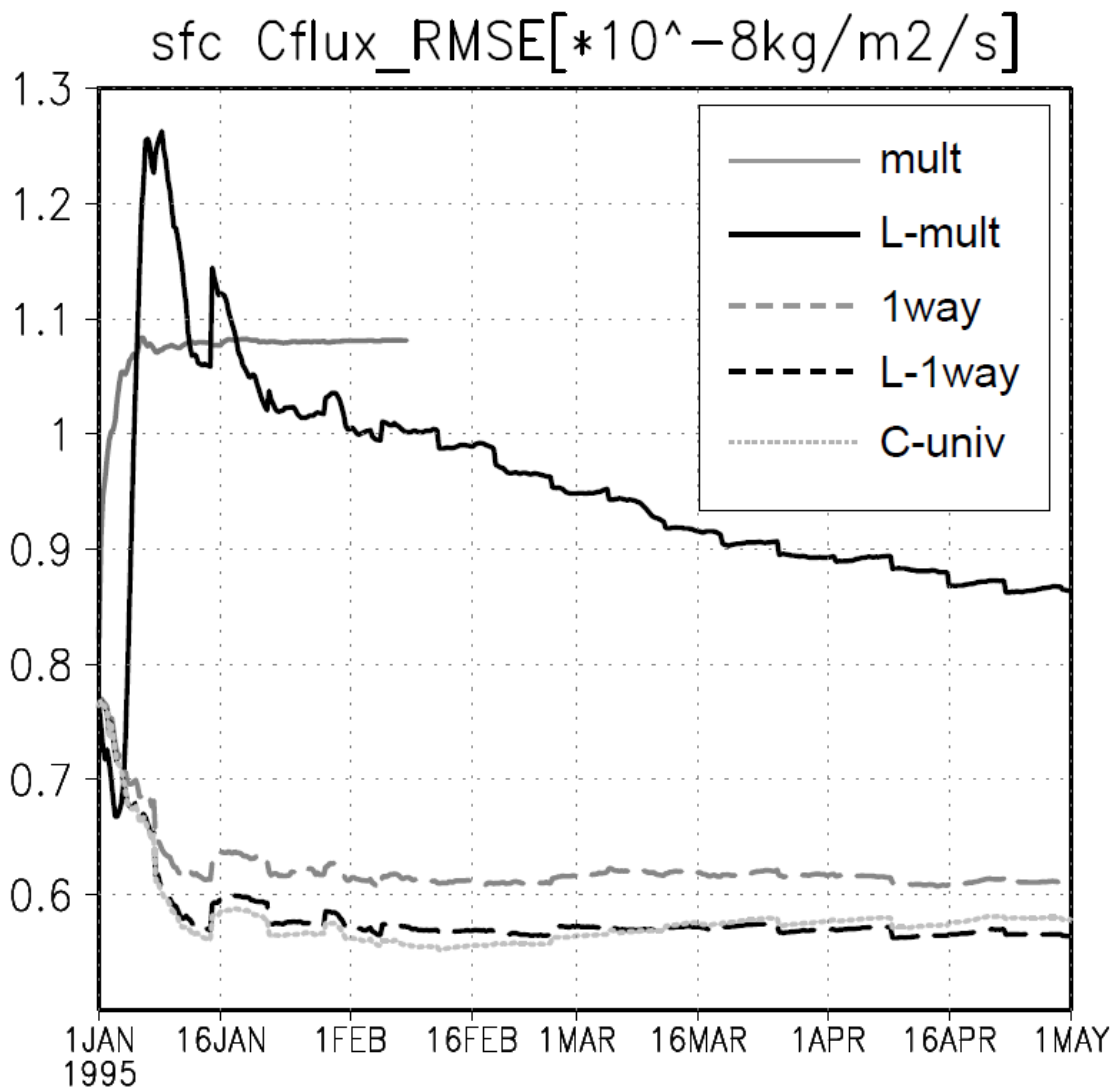
830

831 Figure 3. The simulated observational coverage of (a) meteorological variables (black dots)
 832 and (b) atmospheric CO₂ concentration (gray lines: GOSAT column data, crosses:
 833 continuous *in situ* data, closed circles: weekly flask data).



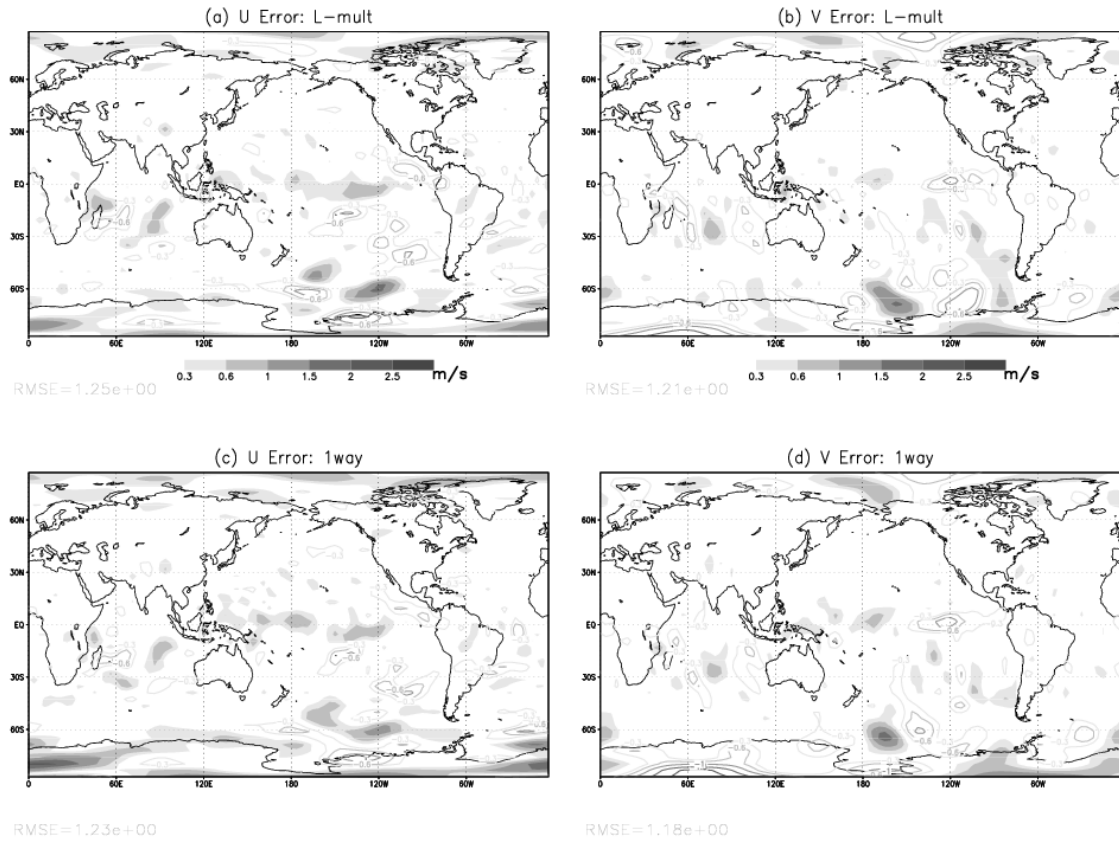
834

835 Figure 4. Time series of global RMS error of (a) U (m/s) and (b) atmospheric CO₂
 836 concentration in the lowest layer (ppmv) for four months of analysis. (solid gray: *mult*,
 837 solid black: *L-mult*, dashed gray: *1way*, dashed black: *L-1way*, dotted light gray: *C-univ*)



838

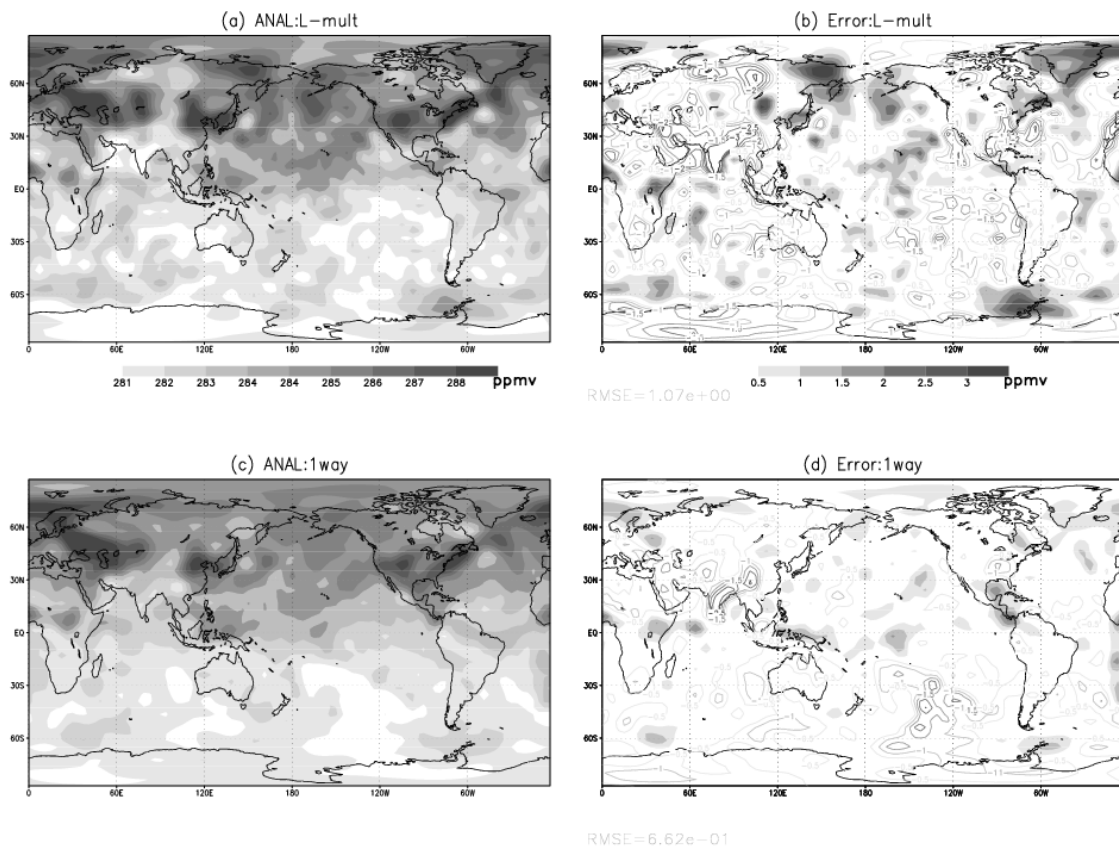
839 Figure 5. Same as Figure 4, except for the surface CO₂ fluxes.



840

841

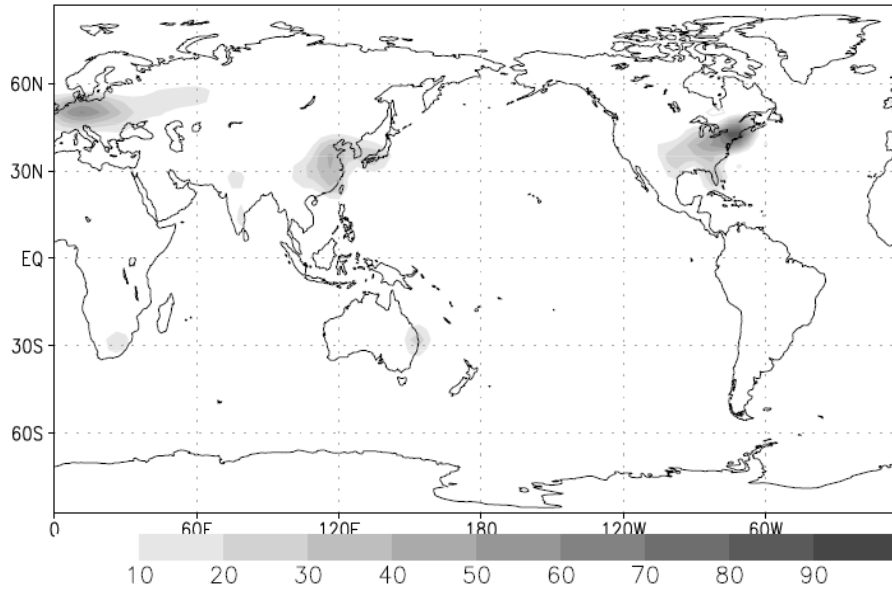
842 Figure 6. Analysis error (unit: m/s) of (a) zonal wind and (b) meridional wind from the
 843 localized multivariate analysis (*L-mult*) for the last three months of data assimilation. (c)
 844 and (d): The same as in (a) and (b) except from *1way*. Shading indicates positive errors and
 845 contours indicate negative errors with the same color scale as the shading.



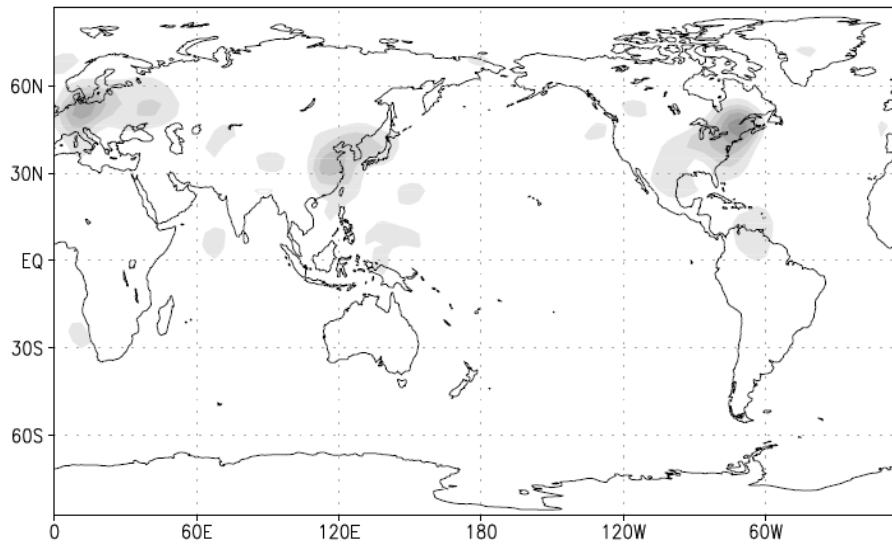
846

847 Figure 7. Analysis (left column) of atmospheric CO₂ concentration (ppmv) in the lowest
 848 layer and its error (right column) after four months of analysis. (a) and (b) results from *L-*
 849 *mult*, (c) and (d) from *1way*. Units are ppmv. Shading indicates positive errors and
 850 contours indicate negative errors with the same color scale as the shading.

(a) True state of CF

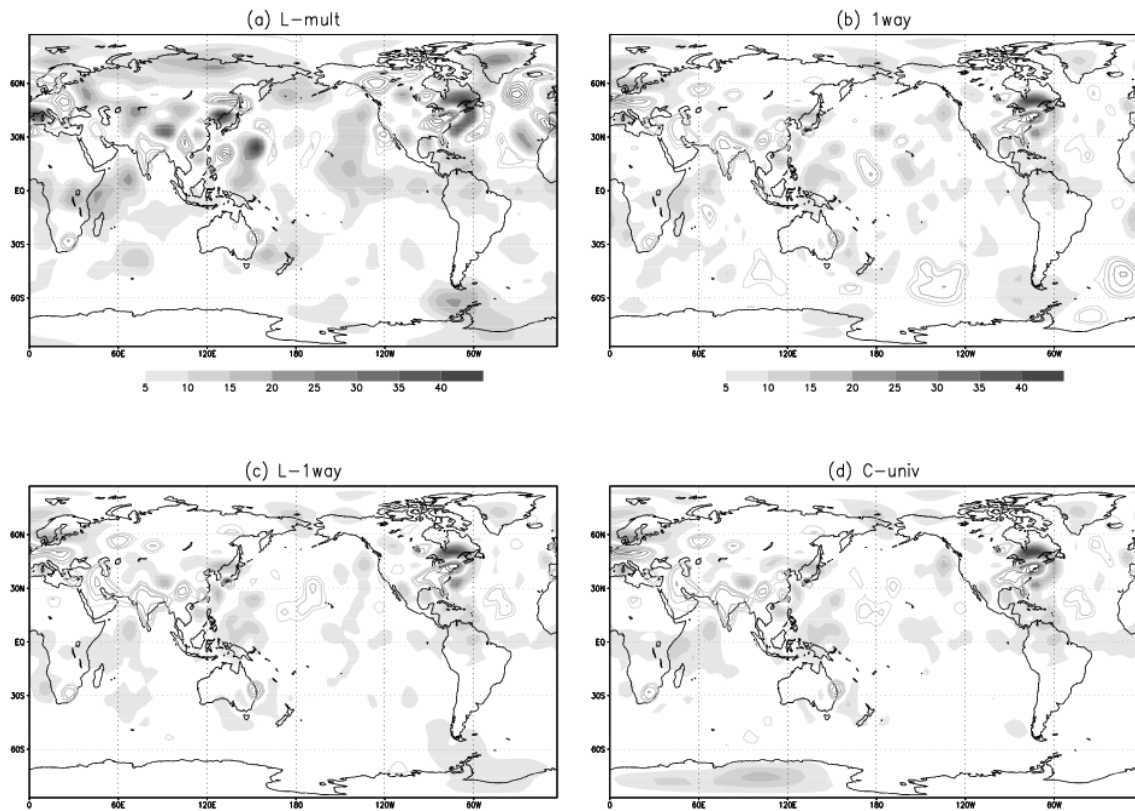


(b) Analysis: L-1way



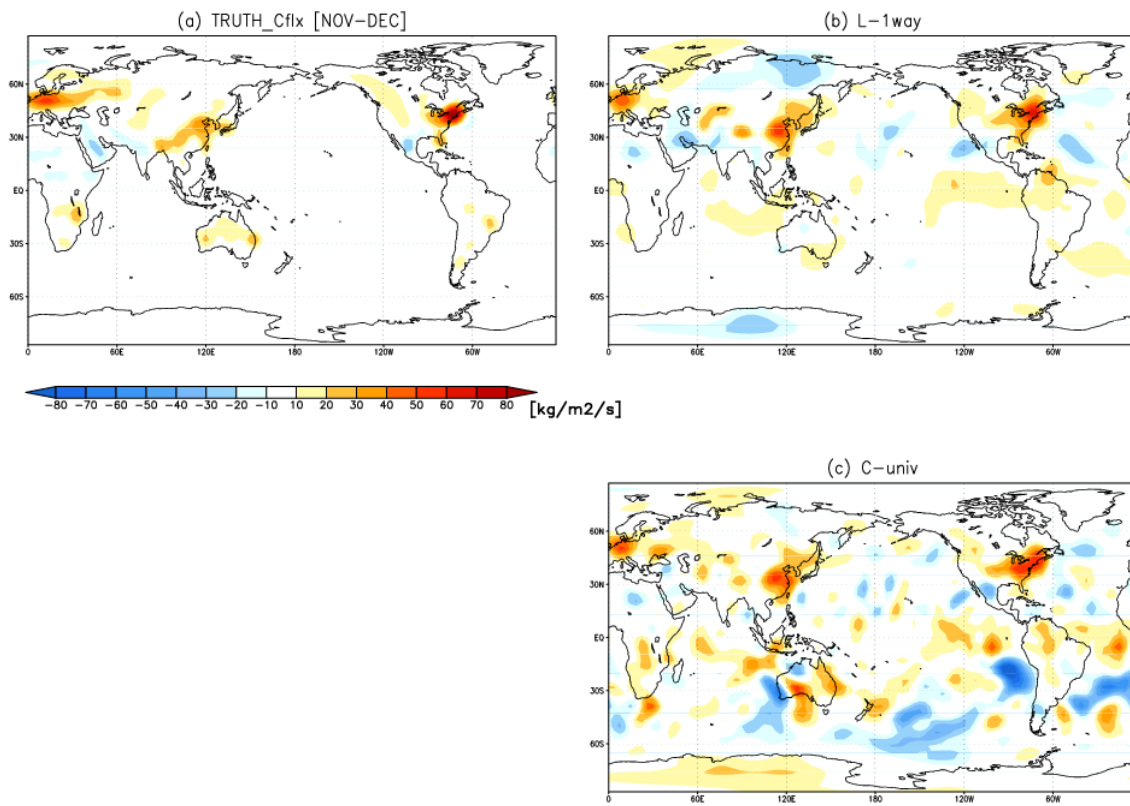
851

852 Figure 8. (a) True state of surface CO₂ fluxes and (b) the analysis after four months of the
853 L-1way (localized 1-way multivariate) data assimilation. (Shading indicates positive errors
854 and contours indicate negative errors with the same contour scale as the shading. Units are
855 10^{-9} kg/m²/s)



856

857 Figure 9. Analysis errors of surface CO₂ fluxes after four months of analysis. (a) results
 858 from *L-mult*, (b) from *1way*, (c) from *L-1way* and (d) from *C-univ*. Units are 10⁻⁹ kg/m²/s.
 859 (Shading indicates positive errors and contours indicate negative errors with the same
 860 contour scale as the shading.)



861

862 Figure 10. (a) True state of surface CO₂ fluxes from a time-varying terrestrial and oceanic
 863 forcing and a fossil fuel emission, and the estimated surface CO₂ fluxes from (b) *L-1way*,
 864 and (c) *C-univ* data assimilation for the last two months (November-December) of one-year
 865 analysis.