

1 **Background error covariance with balance constraints for aerosol species and**
2 **applications in variational data assimilation**

3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33

Zengliang Zang¹, Zilong Hao¹, Yi Li¹, Xiaobin Pan¹, Wei You¹, Zhijin Li² and Dan Chen³

Units: ¹ College of Meteorology and Oceanography, PLA University of Science and Technology,
Nanjing 211101, China;

²Joint Institute For Regional Earth System Science and Engineering, University of California, Los
Angeles, California90095, USA;

³National Center for Atmospheric Research, Boulder, Colorado 80305, USA

June 8, 2016

Corresponding author:
Prof. Zengliang Zang
E-mail: zzlqxy@163.com
Telephone: 86-025-80830400
Fax: 86-025-80830400
Address: No.60, Shuanglong Street, Nanjing 211101, China

34 **Abstract**

35 Balance constraints are important for background error covariance (BEC) in data assimilation
36 to spread information between different variables and produce balance analysis fields. Using
37 statistical regression, we develop a balance constraint for the BEC of aerosol variables and apply it
38 to a three-dimensional variational data assimilation system in the WRF/Chem model. One-month
39 forecasts from the WRF/Chem model are employed for BEC statistics. The cross-correlations
40 between the different species are generally high. The largest correlation occurs between elemental
41 carbon and organic carbon with as large as 0.9. After using the balance constraints, the
42 correlations between the unbalanced variables reduce to less than 0.2. A set of data assimilation
43 and forecasting experiments is performed. In these experiments, surface PM_{2.5} concentrations and
44 speciated concentrations along aircraft flight tracks are assimilated. The analysis increments with
45 the balance constraints show spatial distributions more complex than those without the balance
46 constraints, which is a consequence of the spreading of observation information across variables
47 due to the balance constraints. The forecast skills with the balance constraints show substantial
48 and durable improvements from the 2nd hour to the 16th hour compared with the forecast skills
49 without the balance constraints. The results suggest that the developed balance constraints are
50 important for the aerosol assimilation and forecasting.

51 **Keyword:** aerosol species, WRF/Chem, data assimilation, balance constraint, background error
52 covariance

53

54

55

56

57

58

59

60

61 **1. Introduction**

62 Aerosol data assimilation in chemical transport models has received an increasing amount of
63 attention in recent years as a basic methodology for improving aerosol analysis and forecasting. In
64 a data assimilation system, the background error covariance (BEC) plays a crucial role in the
65 success of an assimilation process. The BEC and the observation error determine analysis
66 increments from the assimilation process (Derber and Bouttier 1999, Chen et al., 2013).

67 However, accurate estimation of the BEC remains difficult due to a lack of information about
68 the true atmospheric states and also due to computational requirement arising from the large
69 dimension of the BEC (typically $10^7 \times 10^7$). For a variational data assimilation system, a few
70 methods have been developed to estimate and simplify the expression of the BEC, such as the
71 analysis of innovations, the NMC (National Meteorological Center) and the ensemble-based
72 (Monte Carlo) methods. The NMC method is extensively used in operational atmospheric and
73 meteorology-chemistry data assimilation systems. It assumes that the forecast errors are
74 approximated by differences between pairs of forecasts that are valid at the same time (Parrish and
75 Derber, 1992). Pagowski et al. (2010) estimated the BEC of PM_{2.5} (particles having an
76 aerodynamic diameter less than 2.5 μm) by calculating the differences between the forecasts of 24
77 and 48 h, and used the estimated BEC in a Grid-point Statistical Interpolation (GSI) system (Wu et
78 al., 2002). Benedetti et al. (2007) estimated the BEC of the sum of the mixing ratios of all aerosol
79 species for an operational analysis and forecast systems at ECMWF (The European Centre for
80 Medium-Range Weather Forecasts). The BEC with multiple species and size bins of aerosols have
81 been calculated and employed in data assimilation. Liu et al. (2011) estimated the BEC with 14
82 aerosol species in the Goddard Chemistry Aerosol Radiation and Transport scheme of the Weather
83 Research and Forecasting/Chemistry (WRF/Chem) model and applied it to the GSI system.
84 Schwartz et al. (2012) increased the number of the species to 15 based on the study of Liu et al.
85 (2011). Li et al. (2013) estimated the BEC for five species derived from the Model for Simulation
86 Aerosol Interactions and Chemistry (MOSAIC) scheme.

87 One important role that the BEC plays in meteorological data assimilation is to spread
88 information between different variables to produce balanced analysis fields, which employ
89 balance constraints to convert original variables into new independent variables. Balance

90 constraints have been employed in atmospheric and oceanic data assimilation, such as geostrophic
91 balance or temperature-salinity balance (Bannister, 2008a, 2008b). To incorporate balance
92 constraints, the model variables are usually transformed to balanced and unbalanced parts. The
93 unbalanced parts as control variables are can be assumed independent, and the balanced parts are
94 constrained by balance constraints (Derber and Bouttier, 1999). Instead of using an empirical
95 function as a balance constraint, balance constraints are also derived using regression techniques
96 (Ricci and Weaver, 2005). Although distinct empirical relations between some variables (such as
97 temperature and humidity) may not exist, the regression equation can also be estimated as balance
98 constraints (Chen et al., 2013).

99 In current aerosol variational data assimilation with multiple variables, balance constraints are
100 not yet incorporated in the BEC. The state variables are assumed to be independent variables
101 without cross-correlation. However, the aerosol species are frequently highly correlated due to
102 their common emission sources and diffusion processes. For example, the correlations in terms of
103 the R-square between elemental carbon and black carbon exceed 0.6 in many locations across Asia
104 and the South Pacific in both urban and suburban locations (Salako et al., 2012), and the
105 correlations between different size bins, such as $PM_{2.5}$ and $PM_{10-2.5}$ (the diameter of particles being
106 between 2.5 and 10 μm), are also generally significant (Sun et al., 2003; Geller et al., 2004). Thus,
107 the cross-correlations between different species or size bins are necessary to produce balanced
108 analysis fields. Cross-correlations spread the observation information from one variable to other
109 variables, which can produce more balanced initial fields. For the data assimilation of the
110 ensemble Kalman filter method, the BEC with balance constraints is assured (Pagowski et al.,
111 2012; Schwartz et al., 2014), although the balance may break down because of localization.

112 Recently, several studies have suggested that the BEC with balanced cross-correlation should
113 be introduced into aerosol variational data assimilation (Kahnert, 2008; Liu et al., 2011; Li et al.,
114 2013; Saide et al., 2013). Kahnert (2008) exhibited cross-correlations of the seventeen aerosol
115 variables from Multiple-scale Atmospheric Transport and Chemistry (MATCH) Model. He found
116 that the statistical cross-correlations between aerosol components are primarily influenced by the
117 interrelations between emissions and by interrelations due to chemical reactions to a much lesser
118 degree. Saide et al., (2012; 2013) incorporated the capacity to add cross-correlations between

119 aerosol size bins in GSI for assimilating observations of aerosol optical depth (AOD) data. The
120 cross-correlations between the two connecting size bins for each species were considered using
121 recursive filters while, the cross-correlation is not considered for the other size bins that are not
122 connecting.

123 In this paper, we explore incorporating cross-correlations between different species in BEC
124 using balance constraints. The balance constraints are established using statistical regression. We
125 apply the BEC with the balance constraints to a data assimilation and forecasting system with the
126 MOSAIC scheme in WRF/Chem. The MOSAIC scheme includes a large number of variables with
127 eight species, and flexibility of eight or four size bins. The scheme of four size bins is used in our
128 studies. The four bins are located between 0.039–0.1 μm , 0.1–1.0 μm , 1.0–2.5 μm , and 2.5–10 μm ,
129 and the total mass of the first three bins are $\text{PM}_{2.5}$. A 3DVAR system for the MOSAIC (4-bin)
130 scheme has been developed by Li et al. (2013). For comparisons, we employ this 3DVAR system
131 with the same model configurations as employed by Li et al. (2013). The data assimilation and
132 forecasting experiments are performed with a focus on assessing the impact of cross-correlations
133 of the BEC on analyses and forecasts.

134 The paper is organized as follows: Section 2 describes the 3DVAR system and the formulation
135 of the BEC. Section 3 describes the WRF/Chem configuration and estimates the correlations
136 among the emissions. The statistical characteristics of the BEC, including the regression
137 coefficient of the cross-correlation, are discussed in Section 4. Using the BEC, experiments of
138 assimilating surface $\text{PM}_{2.5}$ observations and aircraft observations are discussed in Section 5.
139 Shortcomings, conclusions and future perspectives are presented in Section 6.

140 **2. Data assimilation system and BEC**

141 In this section, we present a formulation of the BEC with cross-correlation between different
142 species using a regression technique. Then, the cost function with the new BEC is derived and the
143 calculating factorization of the BEC is described.

144 The control variables of the data assimilation are obtained from the MOSAIC (4-bin) aerosol
145 scheme in the WRF/Chem model (Zaveri et al., 2008). The MOSAIC scheme includes eight
146 aerosol species, that is, elemental carbon or black carbon (EC/BC), organic carbon (OC), nitrate
147 (NO_3), sulfate (SO_4), chloride (Cl), sodium (Na), ammonium (NH_4), and other inorganic mass

148 (OIN). Each species is separated into four bins with different sizes: 0.039–0.1 μm , 0.1–1.0 μm ,
 149 1.0–2.5 μm and 2.5–10 μm . The scheme involves 32 aerosol variables with eight species and four
 150 size bins. These variables cannot be directly introduced as control variables in an assimilation
 151 system in consideration of computational efficiency. The number of variables must be decreased
 152 prior to assimilation. Li et al. (2013) have lumped these variables into five species as control
 153 variables in the 3DVAR system. The five species consist of EC, OC, NO_3 , SO_4 and OTR. Here,
 154 OTR is the sum of Cl, Na, NH_4 and OIN. Note that the data assimilation system aims to assimilate
 155 the observation of $\text{PM}_{2.5}$; only the first three of four size bins are utilized to lump as one control
 156 variable for each species.

157 For a 3DVAR system, the cost function (J), which measures the distance of the state vector to the
 158 background and observations, can be written as follows:

$$159 \quad J(\mathbf{x}) = \frac{1}{2}(\mathbf{x} - \mathbf{x}^b)^T \mathbf{B}^{-1}(\mathbf{x} - \mathbf{x}^b) + \frac{1}{2}(\mathbf{y} - \mathbf{H}\mathbf{x})^T \mathbf{R}^{-1}(\mathbf{y} - \mathbf{H}\mathbf{x}). \quad (1)$$

160 Here, \mathbf{x} is the vector of the state variables, including EC, OC, NO_3 , SO_4 and OTR; \mathbf{x}^b is the
 161 background vector of these five species, which are generated by the MOSAIC scheme; \mathbf{y} is the
 162 observation vector; \mathbf{H} is the observation operator that maps the model space to the observation
 163 space and is assumed to be linear here; \mathbf{R} is the observation error covariance associated with \mathbf{y} ;
 164 and \mathbf{B} is the background error covariance associated with \mathbf{x}^b . Eq. (1) is usually written in the
 165 incremental form

$$166 \quad J(\delta\mathbf{x}) = \frac{1}{2}\delta\mathbf{x}^T \mathbf{B}^{-1}\delta\mathbf{x} + \frac{1}{2}(\mathbf{H}\delta\mathbf{x} - \mathbf{d})^T \mathbf{R}^{-1}(\mathbf{H}\delta\mathbf{x} - \mathbf{d}), \quad (2)$$

167 where $\delta\mathbf{x}$ ($\delta\mathbf{x} = \mathbf{x} - \mathbf{x}^b$) is the incremental state variable. The observation innovation vector is
 168 known as $\mathbf{d} = \mathbf{y} - \mathbf{H}\mathbf{x}^b$. The minimization solution is the analysis increment $\delta\mathbf{x}$, and the final
 169 analysis is $\mathbf{x}^a = \mathbf{x}^b + \delta\mathbf{x}$. This analysis is statistically optimal as a minimum error variance
 170 estimate (e.g., Jazwinski, 1970; Cohn, 1997).

171 In Eq. (1) or Eq. (2), \mathbf{x}^b is a $(N \times m)$ – vector, where N is the number of model grid points,
 172 and m is the number of state variables. \mathbf{B} is a symmetric matrix with a dimension of $(N \times m)^2$.
 173 For a high-resolution model, the number of vector \mathbf{x}^b is on the order of 10^7 . Therefore, the
 174 number of elements in \mathbf{B} is approximately 10^{14} . With this dimension, \mathbf{B} cannot be explicitly
 175 manipulated. To pursue simplifications of \mathbf{B} , we employ the following factorization

176
$$\mathbf{B} = \mathbf{D}\mathbf{C}\mathbf{D}^T, \quad (3)$$

177 where \mathbf{D} and \mathbf{C} are the standard deviation matrix and the correlation matrix, respectively. \mathbf{D}
 178 and \mathbf{C} can be described and separately prescribed after the factorization. \mathbf{D} is a diagonal matrix
 179 whose elements include the standard deviation of all state variables in the three-dimensional grids.
 180 To reduce the computational cost, we use the average value of standard deviations that are at the
 181 same level. Thus, the standard deviation is simplified with vertical levels. \mathbf{C} is a symmetric
 182 matrix, having the form

183
$$\mathbf{C} = \begin{bmatrix} \mathbf{C}_{EC}^{EC} & \mathbf{C}_{EC}^{OC} & \mathbf{C}_{EC}^{NO_3} & \mathbf{C}_{EC}^{SO_4} & \mathbf{C}_{EC}^{OTR} \\ \mathbf{C}_{OC}^{EC} & \mathbf{C}_{OC}^{OC} & \mathbf{C}_{OC}^{NO_3} & \mathbf{C}_{OC}^{SO_4} & \mathbf{C}_{OC}^{OTR} \\ \mathbf{C}_{NO_3}^{EC} & \mathbf{C}_{NO_3}^{OC} & \mathbf{C}_{NO_3}^{NO_3} & \mathbf{C}_{NO_3}^{SO_4} & \mathbf{C}_{NO_3}^{OTR} \\ \mathbf{C}_{SO_4}^{EC} & \mathbf{C}_{SO_4}^{OC} & \mathbf{C}_{SO_4}^{NO_3} & \mathbf{C}_{SO_4}^{SO_4} & \mathbf{C}_{SO_4}^{OTR} \\ \mathbf{C}_{OTR}^{EC} & \mathbf{C}_{OTR}^{OC} & \mathbf{C}_{OTR}^{NO_3} & \mathbf{C}_{OTR}^{SO_4} & \mathbf{C}_{OTR}^{OTR} \end{bmatrix}, \quad (4)$$

184 where \mathbf{C}_{EC} , \mathbf{C}_{OC} , \mathbf{C}_{NO_3} , \mathbf{C}_{SO_4} and \mathbf{C}_{OTR} at diagonal locations are the background error
 185 auto-correlation matrices that are associated with each species. They represent the correlation
 186 among pairs of grid points for one species. Other submatrices represent the correlations between
 187 different species, known as cross-correlations. For example, \mathbf{C}_{OC}^{EC} represents the cross-correlations
 188 between EC and OC, and $\mathbf{C}_{OC}^{EC} = (\mathbf{C}_{EC}^{OC})^T$. In Li et al. (2013), these cross-correlations were
 189 disregarded, that is, the five species are considered independently and \mathbf{C} is thus a block diagonal
 190 matrix.

191 In this study, the cross-correlations between different species are considered by introducing
 192 control variable transforms (Derber and Bouttier, 1999; Barker, 2004; Huang, 2009). We divide
 193 the model aerosol variables into balanced components ($\delta\mathbf{x}_b$) and unbalanced components ($\delta\mathbf{x}_u$):

194
$$\delta\mathbf{x} = \delta\mathbf{x}_b + \delta\mathbf{x}_u. \quad (5)$$

195 Note the EC does not need to be divided. There is not unbalanced component for EC that is
 196 similar to the variable of vorticity in the data assimilation of ECMWF (Derber and Bouttier, 1999),
 197 or the variable of stream function in the data assimilation of MM5 (Barker, 2004). The
 198 transformation from unbalanced variables ($\delta\mathbf{x}_u$) to full variables ($\delta\mathbf{x}$) by the balance operator \mathbf{K}
 199 is given by

200
$$\delta\mathbf{x} = \mathbf{K}\delta\mathbf{x}_u. \quad (6)$$

201 Eq. (6) can be written as

$$202 \begin{bmatrix} \delta EC \\ \delta OC \\ \delta NO_3 \\ \delta SO_4 \\ \delta OTR \end{bmatrix} = \begin{bmatrix} \mathbf{I} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \rho_{21} & \mathbf{I} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \rho_{31} & \rho_{32} & \mathbf{I} & \mathbf{0} & \mathbf{0} \\ \rho_{41} & \rho_{42} & \rho_{43} & \mathbf{I} & \mathbf{0} \\ \rho_{51} & \rho_{52} & \rho_{53} & \rho_{54} & \mathbf{I} \end{bmatrix} \begin{bmatrix} \delta EC \\ \delta OC_u \\ \delta NO_{3u} \\ \delta SO_{4u} \\ \delta OTR_u \end{bmatrix}, \quad (7)$$

203 where ρ_{ij} is the submatrix of \mathbf{K} , which represents the statistical regression coefficients between
 204 the variables i and j (Chen et al., 2013). Note that ρ_{ij} is a diagonal matrix with the dimension of
 205 model grid points. Each model grid point has a regression coefficient. For convenience, we
 206 assumed that the elements of ρ_{ij} is a constant value for all grid points, which are denoted as ρ_{ij}
 207 and are calculated by linear regression with all grid points. For example, ρ_{21} can be obtained
 208 from the regression equation of OC and EC as

$$209 \delta OC = \rho_{21} \delta EC + \varepsilon, \quad (8)$$

210 where ε is the residual. δEC and δOC can be estimated from the forecast differences of 24 h
 211 forecasts and 48 h forecasts, similar to the statistics of the BEC. Eq. (8) contains the slope but no
 212 intercept. The intercept is nearly zero because δEC and δOC represent forecast differences that
 213 can be considered to be zero mean values. After obtaining ρ_{21} , the balanced part (e.g., the value
 214 of the regression prediction) of δOC can be obtained by

$$215 \delta OC_b = \widehat{\delta OC} = \rho_{21} \delta EC. \quad (9)$$

216 Where $\widehat{\delta OC}$ represents the predicted value of Eq. (8), which is equal to the balanced part (δOC_b).
 217 Remove the δOC_b from the full variables to obtain the unbalanced part (δOC_u), that is, ε in Eq.
 218 (8). Thus, the calculation of δOC_u can be written as

$$219 \delta OC_u = \delta OC - \rho_{21} \delta EC. \quad (10)$$

220 Here, δOC_u and δEC are employed as predictors in the next regression equation to obtain
 221 δNO_{3b} . Then, we can obtain the unbalanced parts of the remaining variables, which are defined as
 222 follows:

$$223 \delta NO_{3u} = \delta NO_3 - (\rho_{31} \delta EC + \rho_{32} \delta OC_u), \quad (11)$$

$$224 \delta SO_{4u} = \delta SO_4 - (\rho_{41} \delta EC + \rho_{42} \delta OC_u + \rho_{43} \delta NO_{3u}), \quad (12)$$

$$225 \delta OTR_u = \delta OTR - (\rho_{51} \delta EC + \rho_{52} \delta OC_u + \rho_{53} \delta NO_{3u} + \rho_{54} \delta SO_{4u}), \quad (13)$$

226 The coefficient of determination (R^2) can be employed to measure the fit of these regressions. It
 227 can be expressed as

$$228 \quad R^2 = \frac{SSR}{SST}, \quad (14)$$

229 where SSR and SST are the regression sum of squares and the sum of squares for total,
 230 respectively.

231 These unbalanced parts can be considered to be independent because they are residual and
 232 random. \mathbf{B}_u denotes the unbalanced variables of the BEC and can be factorized as

$$233 \quad \mathbf{B}_u = \mathbf{D}_u \mathbf{C}_u \mathbf{D}_u^T, \quad (15)$$

234 where \mathbf{D}_u and \mathbf{C}_u are the standard deviation matrix and the correlation matrix, respectively. \mathbf{C}_u
 235 should be a block diagonal without cross-correlations as follows:

$$236 \quad \mathbf{C}_u = \begin{bmatrix} \mathbf{C}_{EC} & & & & \\ & \mathbf{C}_{OCu} & & & \\ & & \mathbf{C}_{NO_{3u}} & & \\ & & & \mathbf{C}_{SO_{4u}} & \\ & & & & \mathbf{C}_{OTRu} \end{bmatrix}. \quad (16)$$

237 According the definition of the BEC,

$$238 \quad \mathbf{B} = \langle (\delta \mathbf{x})(\delta \mathbf{x}^T) \rangle. \quad (17)$$

239 And \mathbf{B}_u can be written as

$$240 \quad \mathbf{B}_u = \langle (\delta \mathbf{x}_u)(\delta \mathbf{x}_u^T) \rangle. \quad (18)$$

241 Using Eq. (6), Eq. (17) and Eq. (18), the relationship between \mathbf{B} and \mathbf{B}_u is

$$242 \quad \mathbf{B} = \mathbf{K} \mathbf{B}_u \mathbf{K}^T. \quad (19)$$

243 $\mathbf{B}^{\frac{1}{2}}$ and $\mathbf{B}_u^{\frac{1}{2}}$ are defined as the square root of \mathbf{B} and the square root of \mathbf{B}_u , respectively. Their
 244 transformation is

$$245 \quad \mathbf{B}^{\frac{1}{2}} = \mathbf{K} \mathbf{B}_u^{\frac{1}{2}}. \quad (20)$$

246 Using Eq. (15), Eq. (20) can be written as follows:

$$247 \quad \mathbf{B}^{\frac{1}{2}} = \mathbf{K} \mathbf{D}_u \mathbf{C}_u^{\frac{1}{2}}. \quad (21)$$

248 Generally, a transformed cost function of Eq. (2) is expressed as a function of a preconditioned
 249 state variable:

250
$$J(\delta\mathbf{z}) = \frac{1}{2}\delta\mathbf{z}^T\delta\mathbf{z} + \frac{1}{2}\left(\mathbf{H}\mathbf{B}^{\frac{1}{2}}\delta\mathbf{z} - \mathbf{d}\right)^T \mathbf{R}^{-1}\left(\mathbf{H}\mathbf{B}^{\frac{1}{2}}\delta\mathbf{z} - \mathbf{d}\right). \quad (22)$$

251 Here, $\delta\mathbf{z} = \mathbf{B}^{-\frac{1}{2}}\delta\mathbf{x}$. Using Eq. (21), Eq. (22) can be written as

252
$$J(\delta\mathbf{z}) = \frac{1}{2}\delta\mathbf{z}^T\delta\mathbf{z} + \frac{1}{2}\left(\mathbf{H}\mathbf{K}\mathbf{D}_u\mathbf{C}_u^{\frac{1}{2}}\delta\mathbf{z} - \mathbf{d}\right)^T \mathbf{R}^{-1}\left(\mathbf{H}\mathbf{K}\mathbf{D}_u\mathbf{C}_u^{\frac{1}{2}}\delta\mathbf{z} - \mathbf{d}\right). \quad (23)$$

253 Eq. (23) is the last form of the cost function with the cross-correlation of \mathbf{B} .

254 According to Li et al. (2013), the correlation matrix of the unbalanced parts (\mathbf{C}_u) is factorized as

255
$$\mathbf{C}_u = \mathbf{C}_{ux} \otimes \mathbf{C}_{uy} \otimes \mathbf{C}_{uz}. \quad (24)$$

256 Here, \otimes denotes the Kronecker product, and \mathbf{C}_{ux} , \mathbf{C}_{uy} and \mathbf{C}_{uz} represent the correlation
 257 matrices between gridpoints in the x direction, the y direction, and the z direction, respectively,
 258 with the sizes $n_x \times n_x$, $n_y \times n_y$, and $n_z \times n_z$, respectively. Here, n_x , n_y and n_z represent the
 259 numbers of grid points in the x direction, y direction, and z direction, respectively. This
 260 factorization can decrease the size of the dimension of \mathbf{C}_u . Another desirable property of Eq. (24)
 261 is

262
$$\mathbf{C}_u^{\frac{1}{2}} = \mathbf{C}_{ux}^{\frac{1}{2}} \otimes \mathbf{C}_{uy}^{\frac{1}{2}} \otimes \mathbf{C}_{uz}^{\frac{1}{2}} \quad (25)$$

263 \mathbf{C}_{ux} and \mathbf{C}_{uy} are expressed by Gaussian functions, and \mathbf{C}_{uz} is directly computed from the proxy
 264 data. They will be discussed in Sec 4.2.

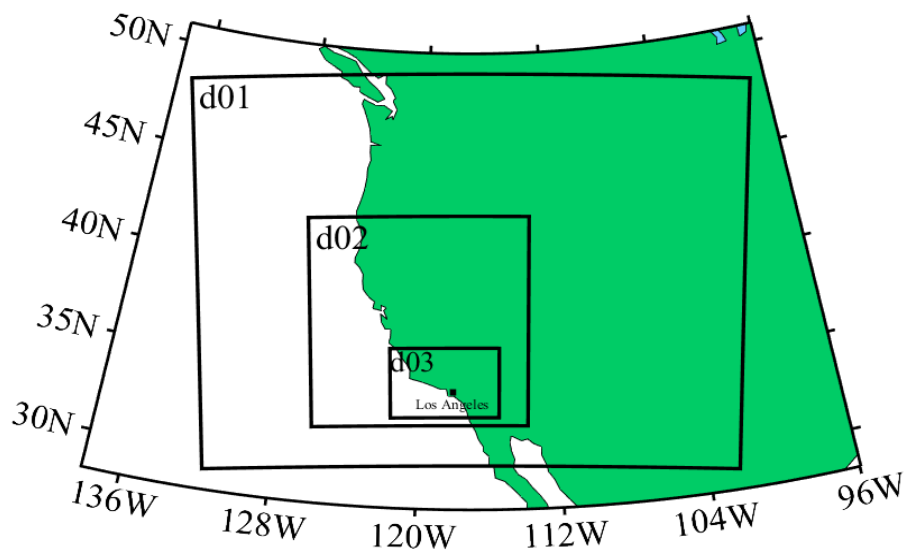
265 3. WRF/Chem configuration and cross-correlations of emission species

266 In this section, we describe the configuration of WRF/Chem, whose forecasting products will
 267 be employed in the following BEC statistics and data assimilation experiments. In addition, the
 268 cross-correlations of emission species from the WRF/Chem emission data are investigated to
 269 understand the cross-correlation between different species of the BEC.

270 3.1 WRF/Chem configuration

271 WRF/Chem (V3.5.1) is employed in our study. This is a fully coupled online model with a
 272 regional meteorological model that is coupled to aerosol and chemistry models (Grell et al., 2005).
 273 The model domain with three spatial domains is shown in Figure 1. The horizontal grid spacing
 274 for these three domains are 36 km (80×60 points), 12 km (97×97 points), and 4 km (144×96
 275 points), respectively. The outer domain spans southern California and the innermost domain

276 encompasses Los Angeles. All domains have 31 vertical levels with the top at 50 hPa. The vertical
277 grid is stretched to place the highest resolution in the lower troposphere. The discussion of the
278 BEC and the emissions presented in this paper will be confined to the innermost domain. The
279 initial meteorology conditions for WRF/Chem are prepared using the North American Regional
280 Reanalysis (NARR) (Mesinger et al. 2006). The meteorology boundary conditions and sea surface
281 temperatures are updated at each initialization. For the forecast running, the initial meteorological
282 conditions are obtained from the NARR data. The initial aerosol conditions are obtained from the
283 former forecast. The emissions are derived from the National Emission Inventory 2005 (NEI'05)
284 for both aerosols and trace gases (Guenther et al., 2006). For more details, the readers are referred
285 to Li et al. (2013).



286

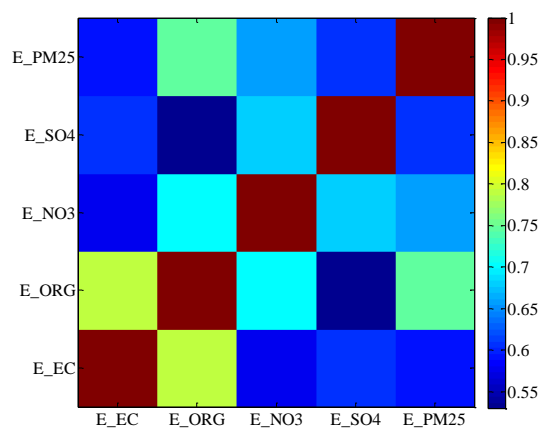
287 Figure 1. Geographical display of the three-nested model domains. The innermost domain covers
288 the Los Angeles basin; the black point denotes the location of Los Angeles.

289 3.2 Cross-correlations of emission species

290 The emission source is necessary for running the WRF/Chem model. It is an important factor
291 for the distribution of the aerosol forecasts. The analysis of the correlations among the emission
292 species can help us to understand the BEC statistics. The emission species is derived from the
293 emission file that is produced by the NEI'05 data for each model domain. Only the emission data
294 for the innermost domain is used to calculate the correlation among the emission species. The
295 emission file contains 37 variables, including gas species and aerosol species. An aerosol species
296 also comprises a nuclei mode and accumulation model species (Peckam et al., 2013). From these

297 aerosol emission species, five lumped aerosol species are calculated, which is consistent with the
 298 variables in the data assimilation. These five lumped species are E_EC (sum of the nuclei mode
 299 and the accumulation mode of elemental carbon PM_{2.5}), E_ORG (sum of the nuclei mode and the
 300 accumulation mode of organic PM_{2.5}), E_NO3 (sum of the nuclei mode and the accumulation
 301 mode of nitrate PM_{2.5}), E_SO4 (sum of the nuclei mode and the accumulation mode of sulfate
 302 PM_{2.5}), and E_PM25 (sum of the nuclei mode and the accumulation mode of unspciated primary
 303 PM_{2.5}).

304 Figure 2 shows the cross-correlations of the five lumped aerosol emission species. All
 305 cross-correlations exceed 0.5. This result reveals that the emission species are correlated, which
 306 may be attributed to the common emission sources and diffusion processes that are controlled by
 307 the same atmospheric circulation. The most significant cross-correlation is between E_EC and
 308 E_ORG with a value of approximately 0.8. This high correlation demonstrates that the emission
 309 distributions of these two species are very similar. Their emissions are primary in urban and
 310 suburban areas with small emissions in rural areas and along roadways (not shown). As shown in
 311 Fig. 2, the lowest cross-correlation is between E_ORG and E_SO4; the latter emissions are
 312 primary in the urban and suburban areas with few emissions in rural areas and roadways (not
 313 shown).



314
 315 Figure 2. Cross-correlations between emission species of E_EC, E_ORG, E_NO3, E_SO4 and
 316 E_PM25. The emission species data are derived from the NEI'05 emissions set for the innermost
 317 domain of the WRF/Chem model

318

319 **4 Balance constraints and BEC statistics**

320 With the configuration of the WRF/Chem model described in Section 3.1, forecasts for one
321 month (from 00UTC of May 15 to 00UTC of June 14, 2010) were performed for the balance
322 constraints and the BEC statistics. Forecast differences between 24 h forecasts and 48 h forecasts
323 are available at 00UTC. Thirty forecast differences are employed as inputs in the NMC method.
324 For this method, 30 forecast differences are sufficient; however, a longer time series may be more
325 beneficial for the BEC statistics (Parrish and Derber, 1992).

326 **4.1 Balance regression statistics**

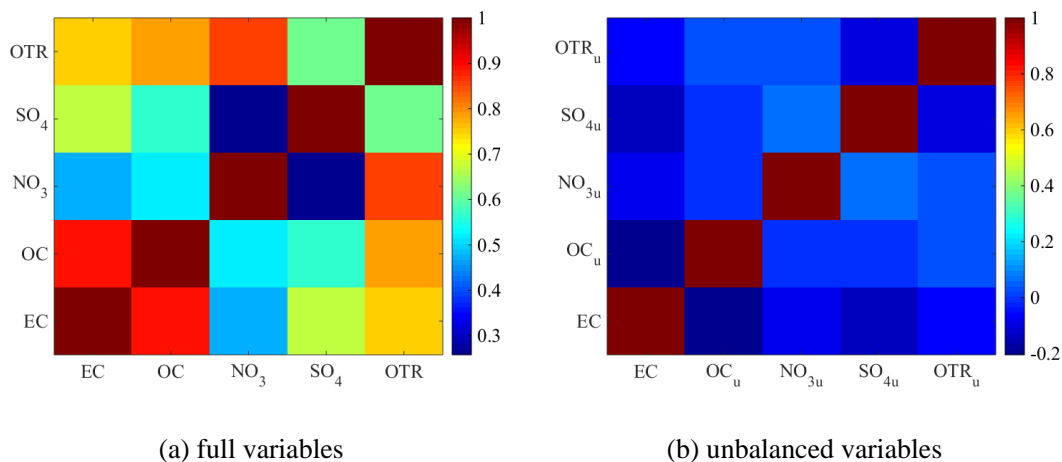
327 Using the 30 forecast differences between 24 h and 48 h forecasts, we can obtain δEC , δOC ,
328 δNO_3 , δSO_4 and δOTR . The size of these variables is $(N \times 30)$, where N is the number of
329 model grid points. We put these data into Eqs. (6-13) to calculate the regression coefficients of ρ_{ij}
330 and the unbalanced parts of the variables. Note the process of calculation should be step by step,
331 since the latter equation will use the unbalanced parts of former equations. Table 1 shows the
332 regression coefficients whose column and row are consistent with $\rho_{i,j}$ in Eq. (7). The last column
333 in Tab. 1 is the coefficient of determination (R^2) of the regression equations. For the regression
334 equation of OC, the regression coefficient is 0.90 and the coefficient of determination of Eq. (7) is
335 0.86, which indicates that EC and OC are highly correlated and their mass concentration scales are
336 approximate. Their correlation is similar to the correlation of the stream function and velocity
337 potential; thus, we set them as the first and second variables in the regression statistics. For the
338 regression equation of NO_3 , the regression coefficients of EC and OC_u are 4.01 and 3.76,
339 respectively, because the mass concentration scale of NO_3 exceeds the mass concentration scales
340 of EC and OC_u . The coefficient of determination is only 0.32, which indicates that the
341 correlations between NO_3 and EC and between NO_3 and OC_u are weak. This result reveals that
342 the forecast errors of NO_3 differ from the forecast errors of EC and OC_u . A possible reason is
343 that NO_3 is the secondary particle that is primarily derived from the transformation of NO_x , but
344 EC and OC_u are derived from direct emissions. Similar to NO_3 , SO_4 is also primarily derived
345 from the transformation of SO_2 and the coefficient of determination for SO_4 is also low. For the
346 last variable OTR, the maximum coefficient of determination is 0.96 because OTR includes some
347 different compositions that are correlated with the first four variables.

348 Table 1 Regression coefficients of balance operator \mathbf{K} and the coefficient of determination
 349 (regression coefficients correspond to ρ_{ij} in Eq. (7))

species	regression coefficient (ρ)					coefficient of determination (R^2)
EC	1					/
OC	0.90	1				0.86
NO ₃	4.01	3.76	1			0.32
SO ₄	1.35	-0.21	-3.15	1		0.48
OTR	2.93	2.35	0.28	0.60	1	0.96

350

351 Figure 3 shows the cross-correlations of the five full variables and the unbalanced variables. In
 352 Fig. 3a, the cross-correlations of the full variables exceed 0.3 and most of them exceed 0.5. In Fig.
 353 3b, however, the cross-correlations of the unbalanced variables are less than 0.2. Some of the
 354 cross-correlations are close to zero, which indicates that these unbalanced variables are
 355 approximatively independent and can be employed as control variables in the data assimilation
 356 system.



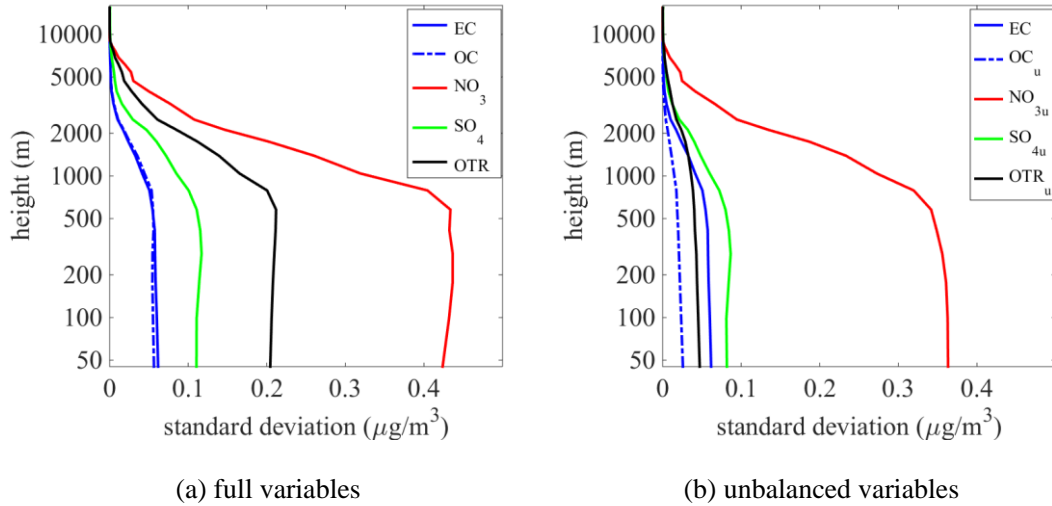
357 Figure 3. Cross-correlations between the five variables of the BEC. These variables are (a) full
 358 variables and (b) unbalanced variables of EC, OC, NO₃, SO₄ and OTR.

359

360 4.2 BEC statistics

361 Using the original full variables and the unbalanced variables obtained by the regression
 362 equations, the BEC statistics are obtained. Figure 4 shows the vertical profiles of the standard

363 deviations of the original \mathbf{D} and the unbalanced \mathbf{D}_u . In Fig. 4a, the original standard deviation of
 364 NO_3 is the largest value, whereas the smallest value is OC, whose profile is close to the profile of
 365 EC. All profiles show a significant decrease at approximately 800 m because the aerosol
 366 particulates are usually limited under the boundary level. In Fig. 4b, all standard deviations
 367 decrease in different degree, with the exception of EC, which remains as the control variable in the
 368 unbalanced BEC statistics. Note that the standard deviation of OTR_u decreases by approximately
 369 80% compared with NO_{3u} , which decreases by approximately 10%. This result is attributed to the
 370 small coefficient of determination for the regression of NO_3 (in Tab. 1), which indicates that a
 371 small portion of NO_3 can be predicted by the regression and a large portion is an unbalanced
 372 component. In contrast with NO_3 , a small portion of OTR is the unbalanced component.



373 Figure 4. Vertical profiles of the standard deviation of the variables. (a) full variables and (b)
 374 unbalanced variables

375

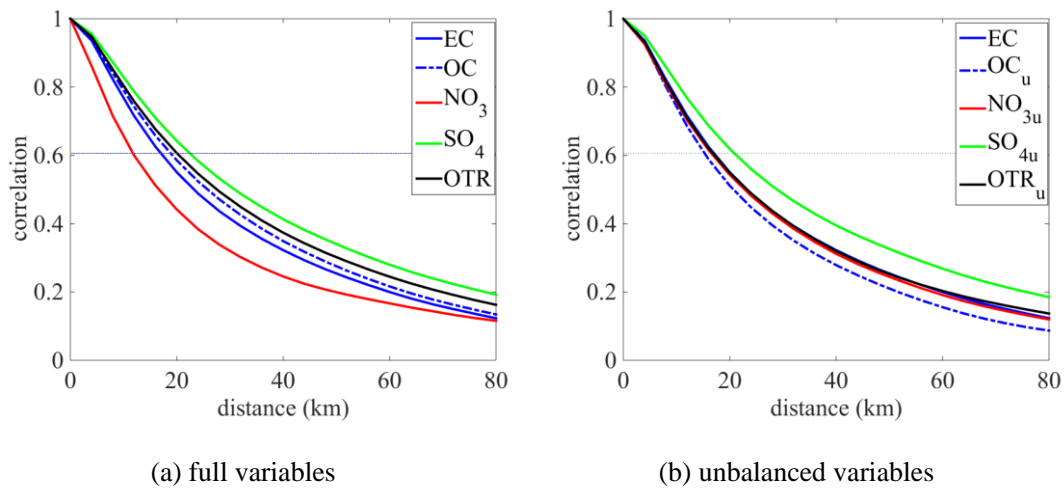
376 For the correlation matrix of \mathbf{C} and \mathbf{C}_u , they are factorized as three independent
 377 one-dimensional correlation matrices in Eq. (24). The horizontal correlation \mathbf{C}_x or \mathbf{C}_y is
 378 approximately expressed by a Gaussian function. The correlation between two points r_1 and r_2

379 can be written as $e^{-\frac{(r_2-r_1)^2}{2L_s^2}}$, where L_s is the horizontal correlation scale and is a constant value

380 for \mathbf{C}_x and \mathbf{C}_y , which are considered to be isotropic (Li et al., 2013). This scale can be estimated

381 by the curve of the horizontal correlations with distances. Figure 5 shows the curves of the

382 horizontal correlations for the five control variables. For the full variables (Fig. 5a), the sharpest
 383 decrease in the curves is observed for NO_3 and the slowest decrease in the curves is observed
 384 for SO_4 . We assume that the decline curve is according to the Gaussian function. Then the
 385 intersection of the decline curve and the line of $e^{-\frac{1}{2}} (\approx 0.61)$ can be approximately as the value
 386 of horizontal correlation scale. The horizontal correlation scales of EC, OC, NO_3 , SO_4 and OTR
 387 are 25 km, 27 km, 20 km, 30 km and 28 km, respectively. For the unbalanced variables (Fig. 5b),
 388 their curves are closer than the curves of the full variables. The correlation scales of EC, OC_u ,
 389 NO_{3u} , SO_{4u} and OTR_u are 25 km, 23 km, 24 km, 28 km and 25 km, respectively. These results
 390 suggest that the unbalanced variables are expressed by some common factors such as EC, OC_u
 391 and NO_{3u} , in the regression equations of Eqs. (10-13), which produces consistent horizontal
 392 correlation scales.

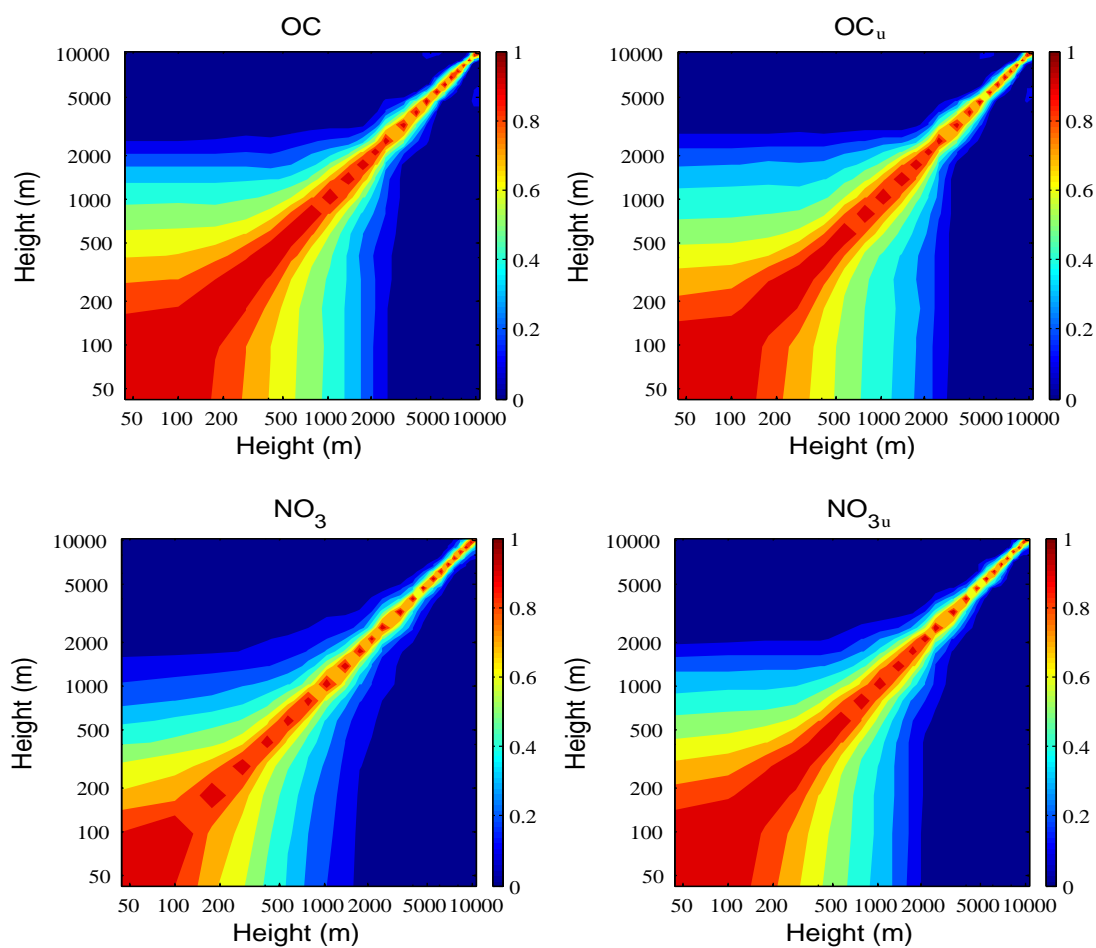


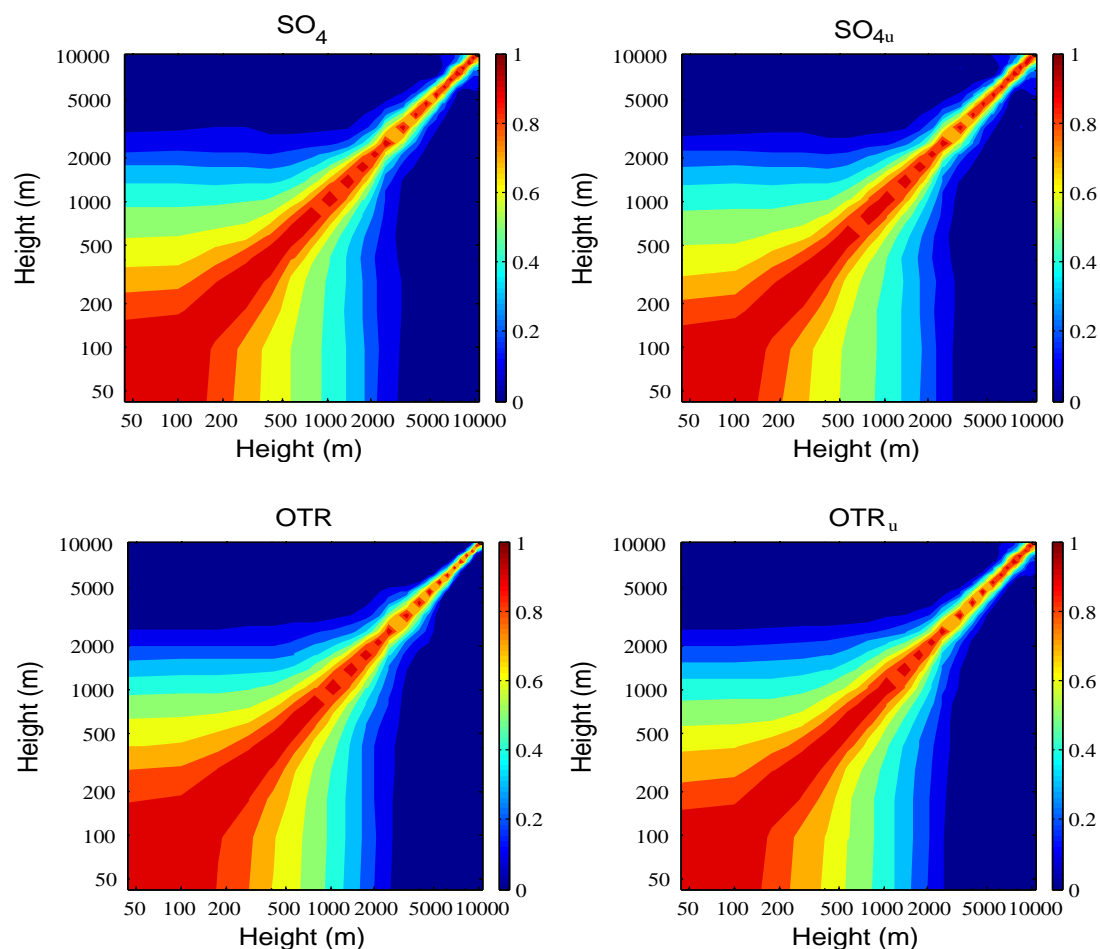
393 Figure 5. Same as Figure 4, with the exception of the horizontal auto-correlation curves of the
 394 variables. The horizontal thin line is the reference line of $e^{-\frac{1}{2}} (\approx 0.61)$ for determining the
 395 horizontal correlation scales.

396

397 For the vertical correlation between \mathbf{C}_z and \mathbf{C}_{uz} , they are directly estimated using the
 398 forecasting differences in the data assimilation system, but not estimated from a approximately
 399 alternative function. Because it is only an $n_z \times n_z$ matrix. Figure 6 shows the vertical correlation
 400 matrices \mathbf{C}_z and \mathbf{C}_{uz} for the full variables (left column) and the unbalanced variables (right
 401 column), respectively. A common feature of both the full variables and the unbalanced variables is

402 the significant correlation between the levels of the boundary layer height, which is consistent
 403 with the profile of the standard deviation in Fig. 4. Some weak adjustments to the correlations
 404 between the full and unbalanced variables are made. For example, the correlation of NO_{3u} is
 405 stronger than the correlation of NO_3 between the boundary layers. Similar with the analysis of
 406 horizontal correlation scale, the vertical correlation scale of NO_{3u} is larger than the vertical
 407 correlation scale of NO_3 . Conversely, the vertical correlation scale of OTR_u is smaller than the
 408 vertical correlation scale of OTR . These results demonstrate that the vertical correlations for the
 409 unbalanced variables are more consistent than the vertical correlations of the full variables, which
 410 is similar to the adjustments to the horizontal correlation scale. Note that the differences of vertical
 411 correlation are slight, compared with the difference of horizontal. The main reason is that the
 412 vertical correlations are generally affected by the atmospheric boundary layer height. Thus, all
 413 vertical correlation decreases rapidly for the levels above the boundary layer height.





414 Figure 6. Vertical correlations of the five variables of the BEC. The left column represents the full
 415 variables, and the right column represents the unbalanced variables.

416

417 5. Application to data assimilation and prediction

418 To exhibit the effect of the balance constraint of the BEC, the data assimilation experiments
 419 and 24-h forecasts for nine cases are run using WRF/Chem model. The surface $PM_{2.5}$ and
 420 aircraft-speciated observations are assimilated using different BEC, and the evaluations are
 421 presented for the data assimilation and subsequent forecasts. Three basic statistical measures
 422 including mean bias (BIAS), root mean square error (RMSE) and correlation coefficient (CORR)
 423 are utilized for the evaluations.

424 5.1 Observation data and experiment scheme

425 Two types of observation data are employed in our experiments. The first type of observation
 426 data consists of hourly surface $PM_{2.5}$ concentrations from the California Air Resources Board
 427 (ARB). There are 42 surface $PM_{2.5}$ monitoring sites existed in the innermost domain of the

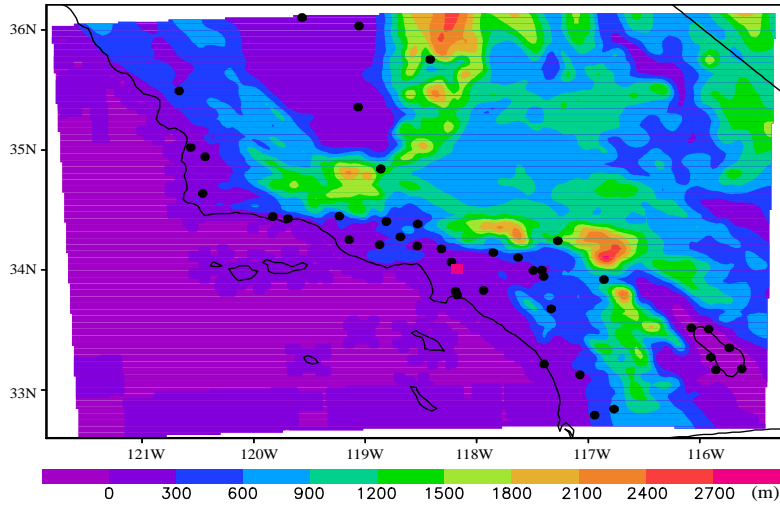
428 WRF/Chem model (Fig. 7). The second type of observation data is the speciated concentration
 429 along the aircraft flight track. The aircraft observations were investigated from the California
 430 Research at the Nexus of Air Quality and Climate Change (CalNex) field campaign in 2010. Nine
 431 flights data around Los Angeles from 15 May to 14 June, 2010 are selected as the cases of data
 432 assimilation. Table 2 shows the start time and end time of each flight. The species of the aircraft
 433 observations include OC, NO₃, SO₄ and NH₄. Note that NH₄ is not a control variable; thus, the
 434 aircraft observation of NH₄ is disregarded in the data assimilation. Because the particle size of the
 435 aircraft observations is less than 1.0 μm, some adjustments to the flight observations are made
 436 according to the ratios between the concentration under 2.5 μm and the concentration under 1.0
 437 μm for each species using model products. With the ratios multiplied by the aircraft observed
 438 concentrations, the speciated concentrations under 2.5 μm can be obtained.

439

440 Table 2 The periods of flight during CalNex 2010 and the initial time of assimilation

Number of cases	Start time of flight	End time of flight	Initial time of assimilation
1	18:00 UTC, May 16	01:42 UTC, May 17	00:00 UTC, May 17
2	17:28 UTC, May 19	00:10 UTC, May 20	18:00 UTC, May 19
3	17:28 UTC, May 21	00:10 UTC, May 21	18:00 UTC, May 21
4	23:08 UTC, May 24	05:23 UTC, May 25	00:00 UTC, May 25
5	01:59 UTC, May 30	07:45 UTC, May 30	06:00 UTC, May 30
6	05:00 UTC, May 31	10:54 UTC, May 31	06:00 UTC, May 31
7	07:59 UTC, June 2	14:09 UTC, June 2	12:00 UTC, June 2
8	07:59 UTC, June 3	14:04 UTC, June 3	12:00 UTC, June 3
9	17:56 UTC, June 14	23:35 UTC, June 14	18:00 UTC, June 14

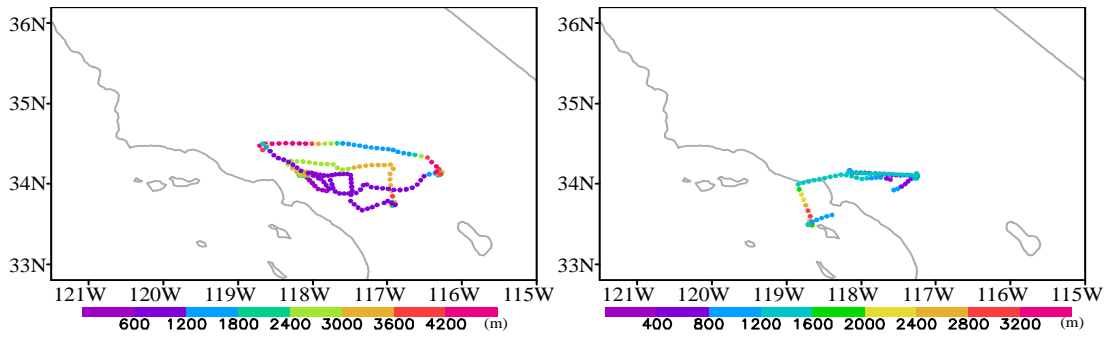
441



442

443 Figure 7. The topography of the innermost domain and the locations of surface monitoring stations
 444 (black dots). The red square is the location of Los Angeles

445

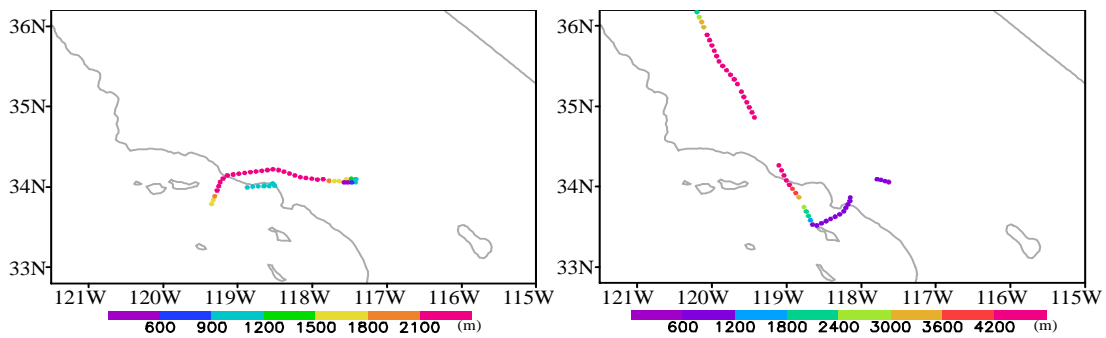


446

447

(a) 00:00 UTC ± 1.5 h, May 17

(b) 18:00 UTC ± 1.5 h, May 19

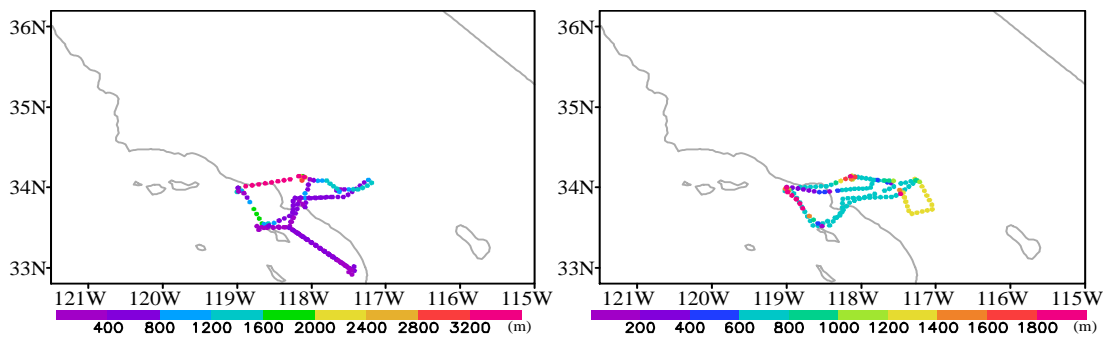


448

449

(c) 18:00 UTC ± 1.5 h, May 21

(d) 00:00 UTC ± 1.5 h, May 25

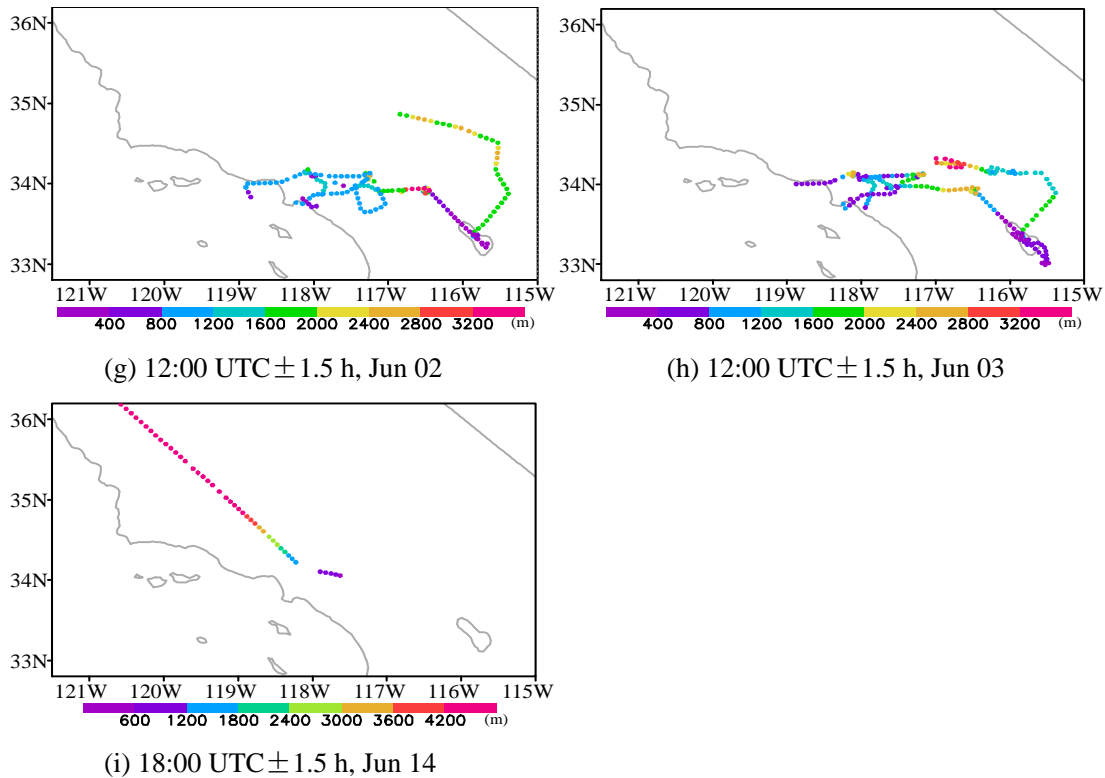


450

451

(e) 06:00 UTC ± 1.5 h, May 30

(f) 06:00 UTC ± 1.5 h, May 31



452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

Figure 8. Aircraft flight tracks during the time window of data assimilation for nine cases. The color of the track indicates the aircraft height.

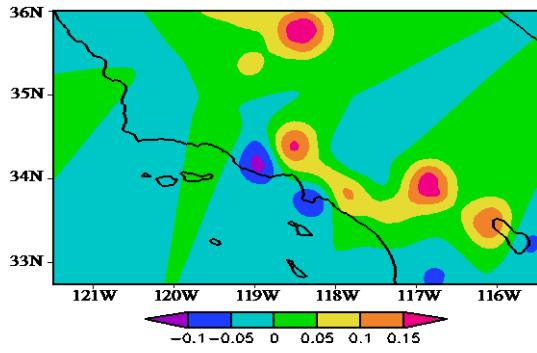
The initial time of data assimilation cases are designed according to the period of flights, showed in Table 2. The time window of assimilation for the flight data is ± 1.5 h, though some flight times do not completely cover the time windows. Figure 8 shows the aircraft tracks during the time window of data assimilation. It is obvious that the aircraft data on May 21, May 25 and June 14 are relative few as the tracks are almost outside of the study domain. For the surface data, it is only the observations at the initial time are assimilated. For each case, three parallel experiments are performed. The first experiment is the control experiment without aerosol data assimilation, which is frequently known as a free run and denoted as Control. The second experiment is a data assimilation experiment that assimilates surface $PM_{2.5}$ and aircraft observations using the full variables without balance constraints; it is denoted as DA-full. The third experiment is also a data assimilation experiment that also assimilates surface $PM_{2.5}$ and aircraft observations, but employs the unbalanced variables as control variables conducted by the balanced constraint; it is denoted as DA-balance. The backgrounds for DA-full and DA-balance are the forecasting results from the previous runs without DA. These previous forecasting results

473 have been obtained when we run the model for the BEC statistics. The observation error is the half
474 of standard deviation of the original background variable, and a vertical profile of observation
475 errors is applied with the average profile of standard deviation of the background variable. In each
476 experiment, a 24-h forecasting is run using the WRF/Chem model with the same configuration
477 described in Section 3.1, and the case on June 3, 2010 is presented in detail as an example.

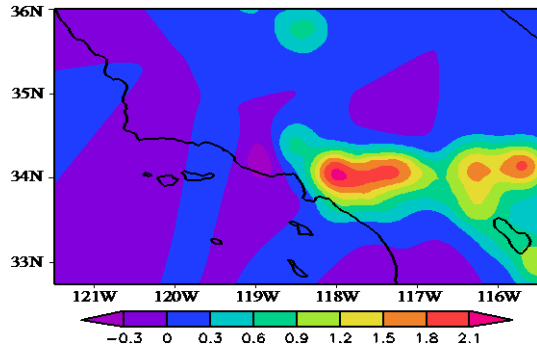
478 **5.2 Increments of data assimilation**

479 Figure 9 shows the horizontal increments of EC, OC, NO₃, SO₄ and OTR at the first model
480 level for the DA-full (left column) and DA-balance experiments (right column) of the case on
481 June 3, 2010. In the DA-full experiment, the increment of EC and OTR (Fig. 9a and 9i) are similar.
482 They are obtained from the surface PM_{2.5} observations because no direct aircraft observations
483 correspond to these two variables. In the DA-balance experiment, significant adjustments are
484 made to the increments of EC (Fig. 9b) under the action of the balance constraints. The
485 observations of OC affect greatly the increments of EC for the high cross-correlation between EC
486 and OC. Thus the increments of EC are similar with the increments of OC. Similarly, significant
487 adjustments are made to the increment of OTR (Fig. 9j), though there are not the species
488 observation of OTR. There are also some slight adjustments for the increments of OC, NO₃ and
489 SO₄ for the crossing spread among species.

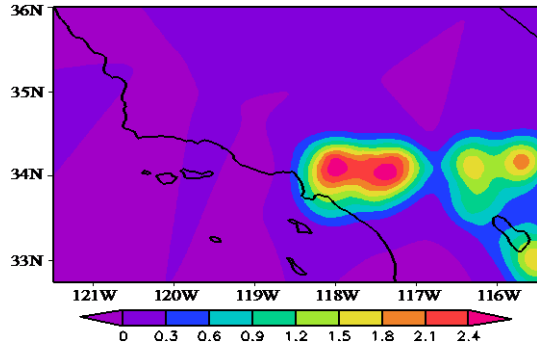
490 Figure 10 shows the vertical increments along 35.0 N for the DA-full and DA-balance
491 experiments. Similar to Fig. 9, the increments of EC and OTR (Fig. 10a and 10i) spread upward
492 from the surface in the DA-full experiment, which are obtained from the surface PM_{2.5}
493 observation. In the DA-balance, the increments of EC and OTR (Fig. 10b and 10j) exhibit
494 observation information from the aircraft height at approximately 500 m, and the value of the
495 increments show significant increases. The distributions of the increments for these five variables
496 in the DA-balance (Fig. 10, right column) generally tend to coincide compared with the
497 distributions of the increments in the DA-full (Fig. 10, left column). The results of the DA-balance
498 are reasonable due to the influence of each other across the balance constraints.



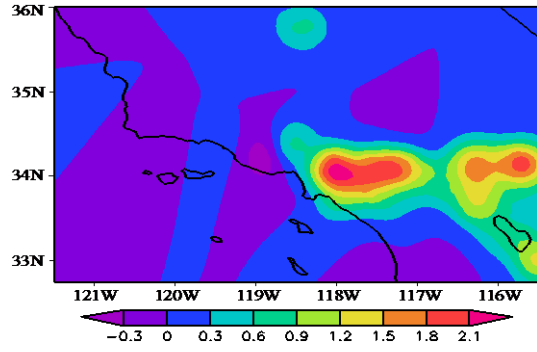
(a) EC in the DA-full



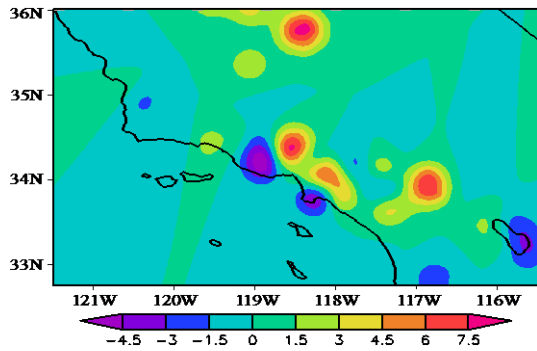
(b) EC in the DA-balance



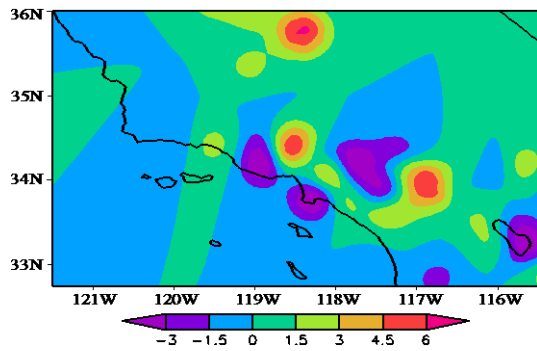
(c) OC in the DA-full



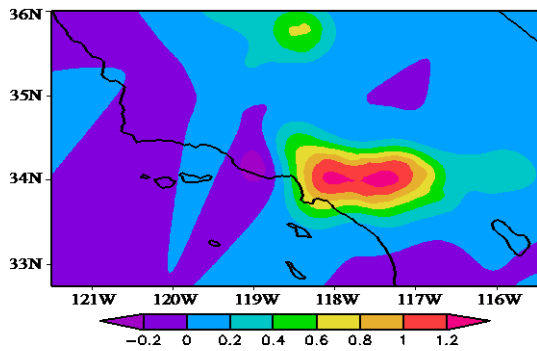
(d) OC in the DA-balance



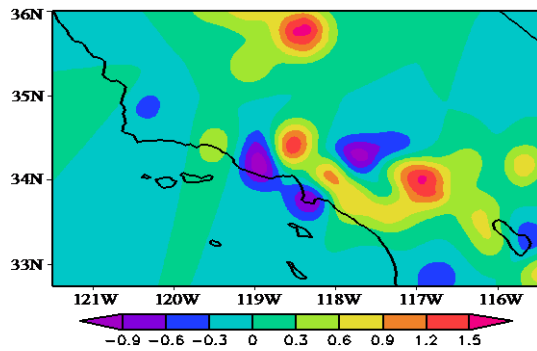
(e) NO₃ in the DA-full



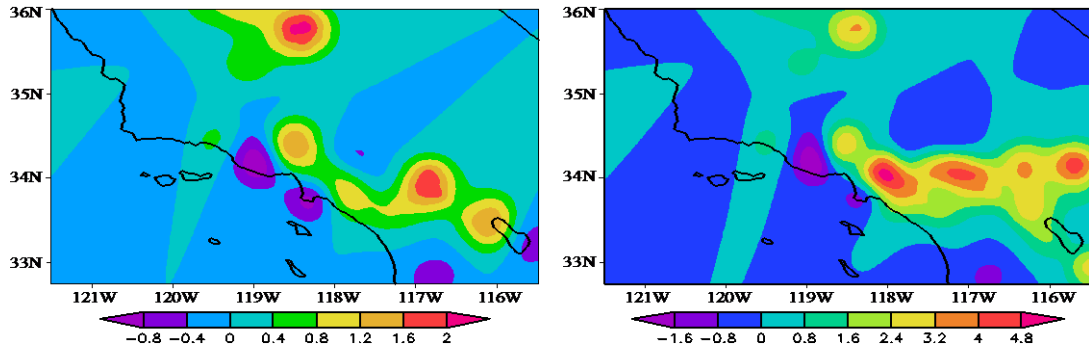
(f) NO₃ in the DA-balance



(g) SO₄ in the DA-full



(h) SO₄ in the DA-balance

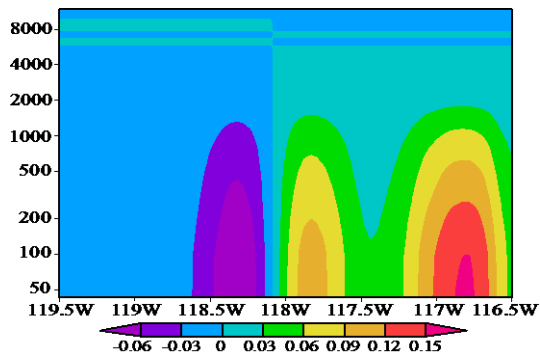


(i) OTR in the DA-full

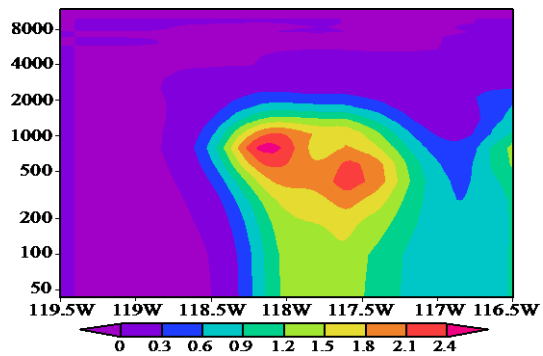
(j) OTR in the DA-balance

499 Figure 9. Surface distributions of increments of the five variables of EC, OC, NO₃, SO₄ and OTR
 500 at 12:00 UTC on June 3, 2010. The left column and right column are from DA-full and
 501 DA-balance, respectively.

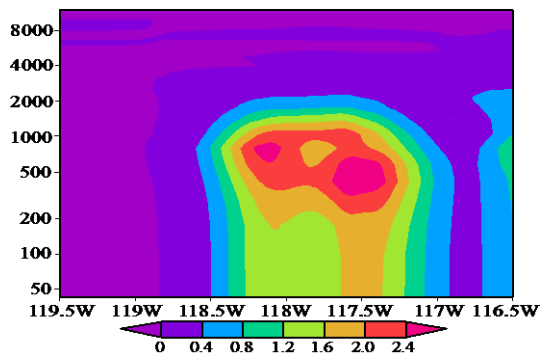
502



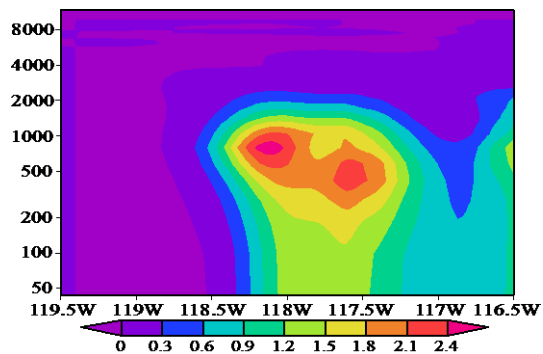
(a) EC in the DA-full



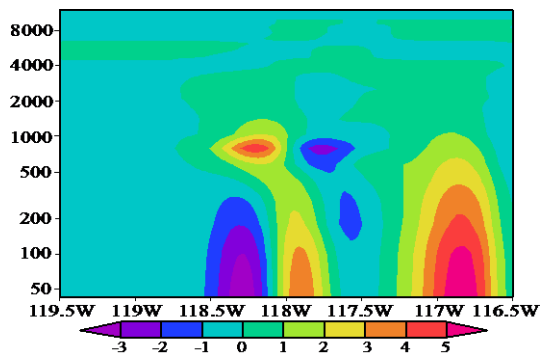
(b) EC in the DA-balance



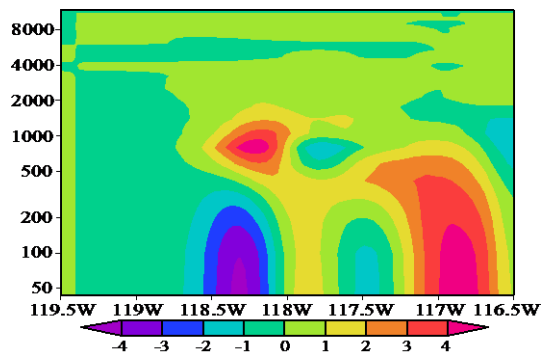
(c) OC in the DA-full



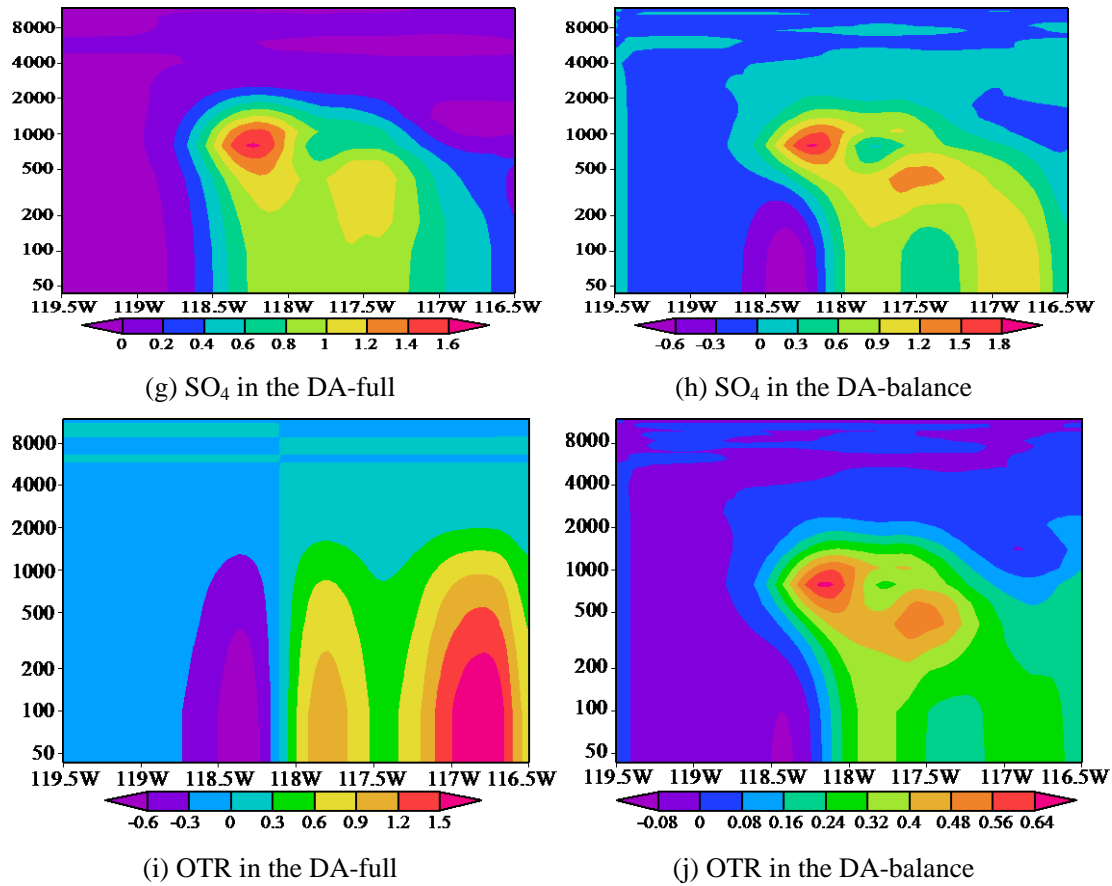
(d) OC in the DA-balance



(e) NO₃ in the DA-full



(f) NO₃ in the DA-balance

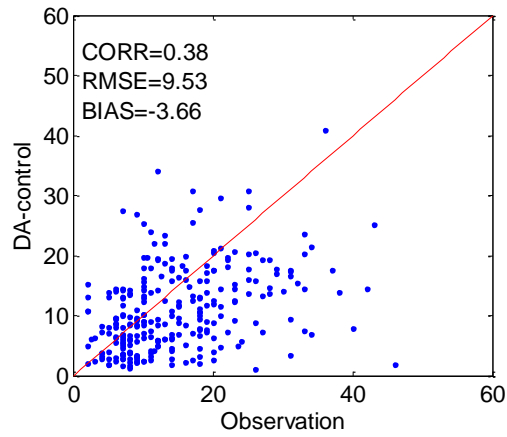


503 Figure 10. Same as Figure 9, with the exception of the vertical sections along 35 N.

504

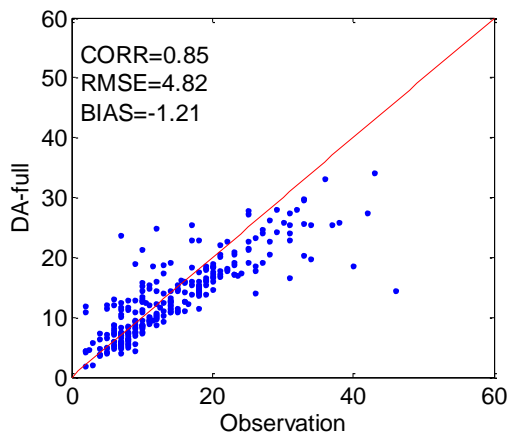
505 **5.3 Evaluation of data assimilation and forecasts**

506 Figure 11 shows the scatter plots of the initial model fields versus the surface observation for all
 507 nine cases. In Fig. 11a, the simulated concentrations of the Control experiment display a
 508 significant underestimation with a BIAS of $-3.66 \mu\text{g}/\text{m}^3$. The mean concentration of Control is
 509 $10.90 \mu\text{g}/\text{m}^3$, about 25.1% lower than observed mean concentrations ($14.56 \mu\text{g}/\text{m}^3$). In the DA-full
 510 and DA-balance experiments, there are remarkable increases for the simulated concentrations, and
 511 the BIASs reduce to as small as -1.21 and $-0.94 \mu\text{g}/\text{m}^3$. The RMSE is $9.53 \mu\text{g}/\text{m}^3$ in the Control
 512 experiment. The RMSE reduces to 4.82 and $4.48 \mu\text{g}/\text{m}^3$ in the DA-full and DA-balance
 513 experiment, respectively. There are also significant improvements for the CORR in the DA-full
 514 and DA-balance experiments, compared with the Control experiment. Furthermore, these three
 515 statistical measures of the DA-balance experiments show some slight improvement, compared
 516 with that of the DA-full experiments. The result demonstrates that more observation information
 517 spread by balance constraints can improve assimilation performance.



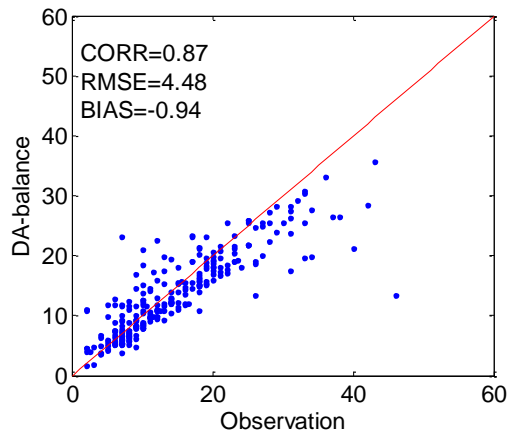
518
519

(a) Control



520
521

(b) DA-full



522
523

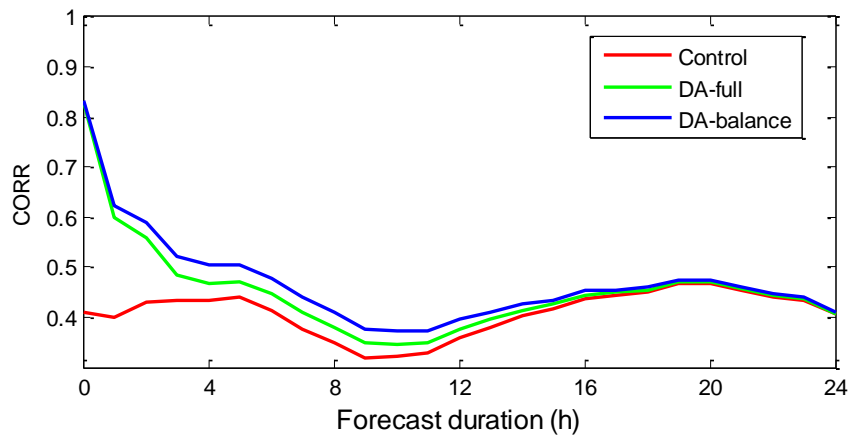
(c) DA-balance

524 Figure 11. Scatter plots of observed concentrations of $PM_{2.5}$ versus simulated $PM_{2.5}$ concentrations
525 of the experiments of (a) Control, (b) DA-full, and (c) DA-balance for all nine cases.

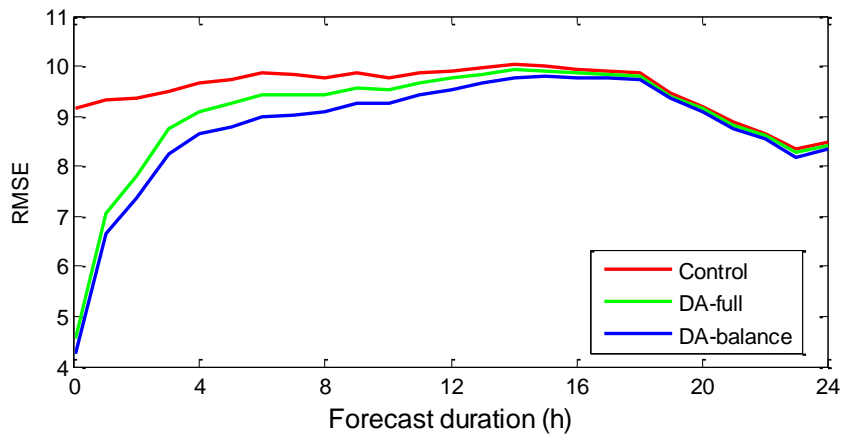
526

527 To evaluate the effects of the data assimilation, the CORR, RMSE and BIAS during the forecast
528 time are calculated for each case, and their averaged results are showed in Figure 12. The CORRs
529 of the DA-balance and DA-full experiments are very close (Fig. 12a). But, the difference increase

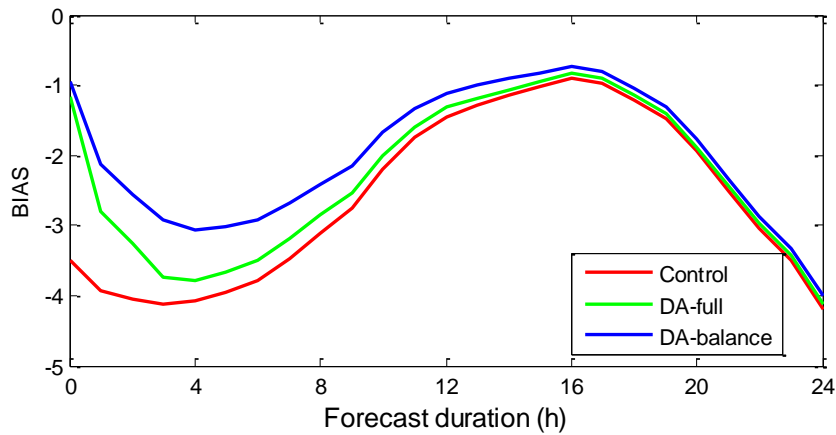
530 after the first hour with a higher CORR in the DA-balance experiment. The CORR of the
531 DA-balance experiment is substantially higher than that of the DA-full experiment from the 2nd
532 hour to the 16th hour. Similar improvements for the RMSE and the BIAS of the DA-balance
533 experiment are observed in Fig. 12b and Fig. 12c, respectively. The improvement for the BIAS in
534 the DA-balance experiment is the most significant among these three statistical measures. The
535 peak value of the improvement for the BIAS (Fig 12c) is at the 4th hour, and the improvement is
536 distinct until the end of forecasts. These improvements indicate that the balance constraint is
537 positive for the subsequent forecasts, which derives from the balanced initial distribution among
538 species.



(a) CORR



(b) RMSE



(c) BIAS

539 Figure 12. The averaged (a) Correlations, (b) root-mean-square errors (RMSE in $\mu\text{g}/\text{m}^3$) and (c)
 540 mean bias (BIAS in $\mu\text{g}/\text{m}^3$) of the $\text{PM}_{2.5}$ concentration forecasts against observations as a function
 541 of forecast duration.

542

543 6. Summary and discussion

544 We examined the BEC in a 3DVAR system, which uses five control variables (EC, OC, NO_3 ,
 545 SO_4 and OTR) that are derived from the MOSAIC aerosol scheme in the WRF/Chem model.
 546 Based on the NMC method, differences within a month-long period between 24- and 48-h
 547 forecasts that are valid at the same time were employed in the estimation and analyses of the BEC.
 548 The background errors of these five control variables are highly correlated. Especially between EC
 549 and OC, their correlation is as large as 0.9.

550 A set of balance constraints was developed using a regression technique and incorporated in
 551 the BEC to account for the large cross correlations. We employ the the balance constraint to
 552 separate the original full variables into balanced and unbalanced parts. The regression technique is
 553 used to express the balanced parts by the unbalanced parts. These unbalanced parts can be
 554 assumed independent. Then, the unbalanced parts are employed as control variables in the BEC
 555 statics. Accordingly, the standard deviations of these unbalanced variables are less than the
 556 standard deviations of the original variables. The horizontal correlation scales of unbalanced
 557 variables are closer than that of full variables on the effect of the balance constraints. And the
 558 vertical correlations of unbalanced variables show similar trend.

559 To evaluate the impact of the balance constraints on the analyses and forecasts, three groups of
 560 experiments, including a control experiment without data assimilation and two data assimilation

561 experiments with and without balance constraints (DA-full and DA-balance), were performed. In
562 the data assimilation experiments, the observations of surface PM_{2.5} concentration and
563 aircraft-speciated concentration of OC, NO₃ and SO₄ were assimilated. The observations of these
564 three variables can spread to the two remaining variables in the increments of the DA-balance,
565 which results in a more complex distribution. The evaluations of CORR, RMSE and BIAS for the
566 initial analysis fields show more improvement in the DA-balance experiments, compared with the
567 DA-full experiments. Though, these improvement are some slight. An important reason is that the
568 surface PM_{2.5} observations are independent from the aircraft observations. If we evaluate the
569 analysis fields by the species observation of aircraft, there may be more significant improvements
570 in the DA-balance experiments.

571 While the improvements increase after the first forecasting hour in the DA-balance
572 experiments, compared with forecasts of the DA-full experiments. The improvements persist to
573 the end of forecasts, and are substantial from the 2nd hour to the 16th hour (Fig. 12). These results
574 suggested that the balance constraints can serve an import role for continually improving the skill
575 of sequent forecasts. Note that some aircraft data are relative few, and some flight tracks are not
576 around Los Angeles in some cases (Fig. 8). If there are more aircraft observations, the
577 improvements of the DA-balance experiments should be more significant and durable.

578 The developed method for incorporating balance constraints in aerosol data assimilation can
579 be employed in other areas or other applications for different aerosol models. For the aerosol
580 variables in different models, some cross-correlations between different species or size bins
581 should exist because their common emissions and diffusion processes are controlled by the same
582 atmospheric circulation. Although these cross-correlations may be stronger than the
583 cross-correlations of atmospheric or oceanic model variables, theoretic balance constraints, such
584 as geostrophic balance or temperature-salinity balance, do not exist. We expected to discover a
585 universal balance constraint that can describe the physical or chemical balanced relationship of
586 aerosol variables, and utilize it in the data assimilation system. In addition, we expected to expand
587 the balance constraint to include gaseous pollutants, such as nitrite (NO₂), sulfur dioxide (SO₂),
588 and (carbon monoxide) CO. These gaseous pollutants are correlated with some aerosol species,
589 such as NO₃, SO₄ and EC, which can improve the data assimilation analysis fields of aerosols by

590 assimilating these gaseous observations. The assimilation of aerosol observations may improve the
591 analysis fields of gaseous pollutants.

592

593 **Code availability**

594 This data assimilation system is established by ourself. The code of this system can be obtained on
595 request from the first author (zzlqxy@163.com).

596

597 **Acknowledgements**

598 This research was supported by the National Natural Science Foundation of China (41275128).
599 We gratefully thank the California Air Resources Board (<http://www.arb.ca.gov/homepage.htm>)
600 and NOAA Earth System Research Laboratory Chemical Sciences Division
601 (<http://esrl.noaa.gov/csd/groups/csd7/measurements/2010calnex/>), for providing the download of
602 surface and aircraft aerosol observations.

603

604 **References**

- 605 Bannister, R.N., 2008a. A review of forecast error covariance statistics in atmospheric variational
606 data assimilation. I: Characteristics and measurements of forecast error covariances. *Quart. J. Roy.
607 Meteor. Soc.*, 134, 1951–1970.
- 608 Bannister, R.N., 2008b. A review of forecast error covariance statistics in atmospheric variational
609 data assimilation. II: Modelling the forecast error covariance statistics. *Quart. J. Roy. Meteor. Soc.*,
610 134, 1971–1996.
- 611 Barker, D.M., Huang, W., Guo, Y.R., Xiao, Q.N., 2004. A Three-Dimensional (3DVAR) data
612 assimilation system for use with MM5: implementation and initial results. *Mon. Weather Rev.* 132,
613 897–914.
- 614 Benedetti, A., Fisher, M., 2007. Background error statistics for aerosols. *Quart. J. Roy. Meteor.
615 Soc.*, 133, 391–405.
- 616 Chen, Y., Rizvi, S., Huang, X., Min, J., Zhang, X., 2013. Balance characteristics of multivariate
617 background error covariances and their impact on analyses and forecasts in tropical and Arctic
618 regions. *Meteorol. Atmos. Phys.*, 121, 79–98.
- 619 Cohn, S.E., 1997: Estimation theory for data assimilation problems: Basic conceptual framework

620 and some open questions. *J. Meteorol. Soc. Jpn.*, 75, 257–288.

621 Derber, J., Bouttier, F., 1999. A reformulation of the background error covariance in the ECMWF
622 global data assimilation system. *Tellus*, 51, 195–221.

623 Geller, M. D., Fine, P. M., and Sioutas, C., 2004. The relationship between real-time and
624 time-integrated coarse (2.5–10m), intermodal (1–2.5m), and fine (<2.5m) particulate matter in the
625 Los Angeles basin. *Journal of the Air & Waste Management Association*, 54(9), 1029–1039.

626 Grell, G.A., Peckham, S.E., Schmitz, R., McKeen, S.A., Frost, G., Skamarock, W.C., Eder, B.,
627 2005. Fully coupled “online” chemistry within the WRF model. *Atmos. Environ.*, 39, 6957–6976,
628 doi:10.1016/j.atmosenv.2005.04.027.

629 Guenther, A., Karl, T., Harley, P., Wiedinmyer, C., Palmer, P.I., Geron, C., 2006. Estimates of
630 global terrestrial isoprene emissions using MEGAN (Model of Emissions of Gases and Aerosols
631 from Nature). *Atmos. Chem. Phys.*, 6, 3181–3210, doi:10.5194/acp-6-3181-2006.

632 Huang, X.Y., Xiao, Q., Barker, D.M., Zhang, X., Michalakes, J., Huang, W., Henderson, T., 2009.
633 Four-dimensional variational data assimilation for WRF: formulation and preliminary results.
634 *Mon. Weather Rev.*, 137, 299–314.

635 Jazwinski, A. H., 1970. *Stochastic processes and filtering theory*, Academic Press, New York, 376
636 pp.

637 Kahnert, M., 2008. Variational data analysis of aerosol species in a regional CTM: background
638 error covariance constraint and aerosol optical observation operators. *Tellus B*, 60(5), 753–770.

639 Li, Z., Zang, Z., Li, Q.B., Chao, Y., Chen, D., Ye, Z., Liu, Y., Liou, K.N., 2013. A
640 three-dimensional variational data assimilation system for multiple aerosol species with
641 WRF/Chem and an application to PM_{2.5} prediction. *Atmos. Chem. Phys.*, 13, 4265–4278.

642 Liu, Z., Liu, Q., Lin, H.C., Schwartz, C.S., Lee, Y.H., Wang, T., 2011. Three-dimensional
643 variational assimilation of MODIS aerosol optical depth: Implementation and application to a dust
644 storm over East Asia. *J. Geophys. Res.*, 116, D23206, doi:10.1029/2011JD016159.

645 Mesinger, F., DiMego, G., Kalnay, E., Shafran, P., Ebisuzaki, W., Jovic, D., Woollen, J., Mitchell,
646 K., Rogers, E., Ek, M., Fan, Y., Grumbine, R., Higgins, W., Li, H., Lin, Y., Manikin, G., Parrish,
647 D., Shi, W., 2006. North American Regional Reanalysis. *B. Am. Meteorol. Soc.*, 87, 343–360.

648 Peckam, S.E., Grell, G.A., McKeen, S.A., Ahmadov, R., 2013. *WRF/Chem Version 3.5 User's*
649 *Guide*. Colorado: NOAA Earth System Research Laboratory.

650 Pagowski, M., Grell, G. A., 2012. Experiments with the assimilation of fine aerosols using an
651 ensemble Kalman filter, *J. Geophys. Res.*, 117, D21302, doi:10.1029/2012JD018333.

652 Pagowski, M., Grell, G.A., McKeen, S.A., Peckham, S.E., Devenyi, D., 2010. Three-dimensional
653 variational data assimilation of ozone and fine particulate matter observations: Some results using
654 the Weather Research and Forecasting–Chemistry model and Grid-point Statistical Interpolation.
655 *Q. J. Roy. Meteorol. Soc.*, 136, 2013–2024, doi:10.1002/qj.700.

656 Parrish, D.F., Derber, J.C., 1992. The national meteorological center spectral statistical
657 interpolation analysis. *Mon. Weather Rev.*, 120, 1747–1763.

658 Ricci, S., Weaver, A.T., 2005. Incorporating State-Dependent Temperature–Salinity Constraints in
659 the Background Error Covariance of Variational Ocean Data Assimilation. *Mon. Weather Rev.*,
660 133, 317–338.

661 Saide, P.E., Carmichael, G.R., Spak, S.N., Minnis, P., Ayers, J.K., 2012. Improving aerosol
662 distributions below clouds by assimilating satellite-retrieved cloud droplet number. *P. Natl. Acad.*
663 *Sci. USA*, 109, 11939–11943, doi:10.1073/pnas.1205877109.

664 Saide, P.E., Carmichael, G.R., Liu, Z., Schwartz, C.S., Lin, H.C., Da Silva, A.M., Hyer, E., 2013.
665 Aerosol optical depth assimilation for a size-resolved sectional model: impacts of observationally
666 constrained, multi-wavelength and fine mode retrievals on regional scale forecasts. *Atmos. Chem.*
667 *Phys.*, 13, 10425–10444, doi:10.5194/acp-13-10425-2013.

668 Salako, G.O., Hopke, P.K., Cohen, D.D., Begum, B.A., Biswas, S.K., Pandit, G.G., Chung, Y.S.,
669 Rahman, S.A., Hamzah, M.S., Davy, P., Markwitz, A., Shagjjamba, D., Lodoysamba, S.,
670 Wimolwattanapun, W., Bunprapob, S., 2012. Exploring the Variation between EC and BC in a
671 Variety of Locations. *Aerosol Air Qual. Res.*, 12, 1–7.

672 Wu, W.S., Purser, R.J., Parrish, D.F., 2002. Three-dimensional variational analysis with spatially
673 inhomogeneous covariances. *Mon. Wea. Rev.*, 130, 2905-2916

674 Schwartz, C.S., Liu, Z., Lin, H.-C., McKeen, S.A., 2012. Simultaneous three-dimensional
675 variational assimilation of surface fine particulate matter and MODIS aerosol optical depth. *J.*
676 *Geophys. Res.*, 117, D13202, doi:10.1029/2011JD017383.

677 Schwartz, C. S., Liu, Z., Lin, H.-C., Cetola, J. D., 2014, Assimilating aerosol observations with a
678 “hybrid” variational-ensemble data assimilation system, *J. Geophys. Res. Atmos.*, 119, 4043–4069,
679 doi:10.1002/ 2013JD020937.

680 Sun C.-H., Lin Y.-C., Wang C.-S., 2003. 7. Relationships among Particle Fractions of Urban and
681 Non-urban Aerosols. *Aerosol and Air Quality Research*, 3(1), 7-15.

682 Zaveri, R.A., Easter, R.C., Fast, J.D., Peters. L K., 2008. Model for Simulating Aerosol
683 Interactions and Chemistry(MOSAIC). *J. Geophys. Res.*, 113, D13204,
684 doi:10.1029/2007JD008782.

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711 Table 1 Regression coefficients of balance operator K and the coefficient of determination712 (regression coefficients correspond to ρ_{ij} in Eq. (7))

species	regression coefficient (ρ)					coefficient of determination (R^2)
EC	1					/
OC	0.90	1				0.86
NO ₃	4.01	3.76	1			0.32
SO ₄	1.35	-0.21	-3.15	1		0.48
OTR	2.93	2.35	0.28	0.60	1	0.96

713

714

715

716

717 Table 2 The periods of flight during CalNex 2010 and the initial time of assimilation

Number of cases	Start time of flight	End time of flight	Initial time of assimilation
1	18:00 UTC, May 16	01:42 UTC, May 17	00:00 UTC, May 17
2	17:28 UTC, May 19	00:10 UTC, May 20	18:00 UTC, May 19
3	17:28 UTC, May 21	00:10 UTC, May 21	18:00 UTC, May 21
4	23:08 UTC, May 24	05:23 UTC, May 25	00:00 UTC, May 25
5	01:59 UTC, May 30	07:45 UTC, May 30	06:00 UTC, May 30
6	05:00 UTC, May 31	10:54 UTC, May 31	06:00 UTC, May 31
7	07:59 UTC, June 2	14:09 UTC, June 2	12:00 UTC, June 2
8	07:59 UTC, June 3	14:041 UTC, June 3	12:00 UTC, June 3
9	17:56 UTC, June 14	23:35 UTC, June 14	18:00 UTC, June 14

718

719

720

721

722

723 Figure 1 Geographical display of the three-nested model domains. The innermost domain covers
724 the Los Angeles basin; the black point denotes the location of Los Angeles.

725

726 Figure 2 Cross-correlations between emission species of E_EC, E_ORG, E_NO3, E_SO4 and
727 E_PM25. The emission species data are derived from the NEI'05 emissions set for the innermost
728 domain of the WRF/Chem model

729

730 Figure 3 Cross-correlations between the five variables of the BEC. These variables are (a) full
731 variables and (b) unbalanced variables of EC, OC, NO₃, SO₄ and OTR.

732

733 Figure 4 Vertical profiles of the standard deviation of the variables. (a) full variables and (b)
734 unbalanced variables

735

736 Figure 5 Same as Figure 4, with the exception of the horizontal auto-correlation curves of the
737 variables. The horizontal thin line is the reference line of $e^{-\frac{1}{2}}$ (≈ 0.61) for determining the
738 horizontal correlation scales.

739

740

741 Figure 6 Vertical correlations of the five variables of the BEC. The left column represents the full
742 variables, and the right column represents the unbalanced variables.

743

744 Figure 7 The topography of the innermost domain and the locations of surface monitoring stations
745 (black dots). The red square is the location of Los Angeles

746

747 Figure 8 Aircraft flight tracks during the time window of data assimilation for nine cases. The
748 color of the track indicates the aircraft height.

749

750 Figure 9 Surface distributions of increments of the five variables of EC, OC, NO₃, SO₄ and OTR
751 at 12:00 UTC on June 3, 2010. The left column and right column are from DA-full and
752 DA-balance, respectively.

753

754 Figure 10 Same as Figure 9, with the exception of the vertical sections along 35 N.

755

756 Figure 11 Scatter plots of observed concentrations of $PM_{2.5}$ versus simulated $PM_{2.5}$ concentrations
757 of the experiments of (a) Control, (b) DA-full, and (c) DA-balance for all nine cases.

758

759 Figure 12 The averaged (a) Correlations, (b) root-mean-square errors (RMSE in $\mu\text{g}/\text{m}^3$) and (c)
760 mean bias (BIAS in $\mu\text{g}/\text{m}^3$) of the $PM_{2.5}$ concentration forecasts against observations as a function
761 of forecast duration.

762