**Referee #1**

Received and published: 12 January 2016

This paper presents the treatment of aerosol background error covariance with balance constraints. Overall, the paper is well written. The presented result shows that the method would improve the chemical data assimilation performance.

<span style="color:red">We thank Reviewer #1 for thoroughly reviewing the manuscript, valuable comments and constructive suggestions. We have carefully addressed all Reviewer's comments and suggestions. We also respond point by point to the reviewer's comments as listed below.</span>

One major concern is that the numerical experiments were only based on a 24-hour forecasting. Since the atmospheric chemistry and meteorological conditions vary day to day. It is highly suggested that the authors extend the experiments to a longer time period. The test period is coincident with CalNex field campaign. So it is not difficult to find more observations for such testing.

<span style="color:red">The period of the effect of data assimilation is generally less than 24 hours (Fig. 12). A longer time forecast of the DA experiment should be very close to the experiment without DA. To demonstrate the robustness of our DA system, we conducted nine cases with a group of 24-h forecasts for each case. For the flight events are discontinuous, we ran these cases with different initial time according the flight processes. The details of these nine cases are in Table 2 (Page 19) and Figure 8 (Page 21). And the results are showed in Figure 11 (Page 26) and Figure 12 (Page 28).</span>

Cross-correlations can be between different species/bins or different grid points. The authors often use "cross-correlation" without specifying what they mean. It is helpful to be unambiguous. For instance, in abstract (line 7 on Page 10054), "cross-correlation" probably refers to the correlation between different species.

<span style="color:red">Thanks, the specification for cross-correlation has been revised in this sentence and some other sentences (line 40, line 123, line 141, line 191, line 268, line 269 and line 580).</span>

The description of model configuration is lacking. Although readers are referred to Li et al. (2013) for details, some basic information should be provided directly. For instance, the mapping projection used for the horizontal coordinates and the extensions of the vertical levels are better given in the paper.

<span style="color:red">We have added some description of the model configuration in the revised manuscript (lines 273-277).</span>

The calculation of the cross-correlation of emission species is not clear. Is it based on 15 May-14 June, 2010 emissions over domain d03?

The emission species are referred to those RADM2 species that produced by NEI'05 data. We calculated the correlation between any two horizontal fields of emission species over domain d03. It is indicated at line 294.

10054, line 14, "are more coincident" -> have better agreement
Thanks, the sentence has been revised (line 47-49).

10054, line 23, "meteorology-chemistry models" -> Chemical transport models
The "meteorology-chemistry models" has been revised as "Chemical transport models" (lines 62).

10055, line 4, "difficult dealt due to ...": Remove "dealt".
The "dealt" has been removed.

10055, line 26, "balance analysis fields": balanced analysis fields?
Corrected.

10056, line 16: PM2.5 is part of PM10 and PM1 is part of PM2.5. So they do not represent different size bins.
This sentence has been revised, and we cited two new papers about the relationship of PM2.5 and PM10-2.5 (lines 109-111).

10056, line 19: The spread of observation impact is not necessarily "enhanced".
Corrected. We have changed to "produce more balanced initial fields" (line 109).

10057, line 4: It is not clear what "the species that are not ADJACENT" means here.
"that are not ADJACENT" means that are not the connecting. This sentence has been revised (line 121-122).

10057, line 12, ".. has been ESTIMATED ...": Developed or applied?
The "estimated" has been changed to "developed" in the revised manuscript.

10058, Eq(1): It is better to have the LHS written as J(x) and x should be in bold font.
"$J(\delta x)$" has been changed to "$J(x)$", and $x$ is in bold font in the revised manuscript.

10059, line 2, "d = y - Hx ": d= y-H x^b
Corrected.

10059, lines 6-7: This seems to neglect the fact that there are multiple variables at each grid point.

Corrected. The sentence has been change to: For a high-resolution model, the number of vector $x^b$ is on the order of $10^7$. Therefore, the number of elements in **B** is approximately $10^{14}$ (line 174).

10059, line 18, "which represent ...": Separate the run-together sentence. They represent the correlation among pairs of grid points for one species.
Thanks. The sentence has been separated (lines 185-186).

10063, line 16, "cross-correlations between emissions": Change to "cross-correlations of emission species", to be consistent with the title of Section 3.2.
The sentence has been revised (Line 265).

10063, line 19, "the cross-correlations of aerosol emissions from ..." -> the cross-correlations of aerosol emission species from ...
The sentence has been revised (Line 268).

10064, line 3, "...that is coupled to aerosol and chemistry domains" -> ...that is coupled to aerosol and chemistry models
The sentence has been revised (Line 272).

10064, lines 16-17, "Emission files are .... of the aerosol forecasts": What does "a primary factor for the distribution of the aerosol forecasts" mean? In addition, it is a run-together sentence that needs to be rewritten.
"A primary factor" means the emission file. This sentence has been revised as: The emission files are necessary for running the WRF/Chem model. It is an important factor for the distribution of the aerosol forecasts. (lines 290-291)

10065, line 5: "With the exception of the auto-correlation in the diagonal line" is redundant.
The sentence has been removed.

10066, line 25: Please spell out "DA" as "data assimilation" since DA is not previously defined yet.
Corrected.

10067, line 3: Figure 2 -> Figure 4.
Corrected.

10067, line 4: Fig. 2a -> Fig. 4a.
Corrected.

10067, line 8: Fig. 2b -> Fig. 4b.
Corrected.

10067, line 8, "all standard deviations significantly decrease": The decrease of NO3 is not significant.
The sentence has been change to: "all standard deviations decrease in different degrees" (Line 367).

10069, line 8: There are many other flights available. Why was this flight chosen over all the others? More description on the flight observations is needed as well.
We have added all flight events that performed during May 15 to 00UTC of June 14, 2010. We chose the case of June 3 for the aircraft observations are enough and the aircraft tracks around Los Angeles, the center of model domain. For the other cases, there are not enough aircraft observations during the assimilation time windows (±1.5 hour of initial time), or the flight tracks are relative few and not around Los Angeles (Fig. 8). But, to demonstrate the robustness of our DA system, we run all cases and calculate the average improvements.

10069, line 25: WRf/Chem -> WRF-Chem. Note that the WRF/Chem is better changed to WRF-Chem in the entire paper.
All "WRF-Chem" have been change to "WRF/Chem".

10071, line 17: Figure 11 does NOT show "scatter plots".
Corrected, thank you.

10071, line 19: Fig.1a -> Fig. 11a
Corrected, thank you.

10072, lines 9-10: It is not true that the data assimilation has the hypothesis of the independent control variables. The independence of control variables merely helps to simplify the background error covariance matrix.
We agree. This sentence has been removed.

10077, Table 1: Please add the name of the species to the table.
Corrected, thank you (Page 33).

10079, Figure 2: Why aren't the cells identical in shape? It applies to Figure 3 too.
Figure 2 and Figure 3 have been plotted in the same shape.

10081, Figure 4, "Same as Fig. 3": Figure 4 is quite different from Figure 3.

Corrected, thank you (line 373-374).

**Referee #2**

This paper discusses implementation of cross-correlations between aerosol variables in a variational data assimilation (DA) system via balance constraint. The authors describe their methodology and then apply their new developments for a single 24-hr forecast over Southern California. Results suggest that incorporating cross-correlations within the DA system was beneficial, especially for 3- to 18-hr forecasts.

This paper is generally interesting and good, although there are some shortcomings that I believe should be addressed before publication. My biggest concerns regard that only one forecast was produced and lack of discussion about other methods of dealing with cross-correlations for aerosols, such as ensemble-based DA methods. Additionally, there are many areas of text that I believe require some clarification.

Thank you very much for your careful review and constructive suggestions. Please find below our detailed responses to all questions and comments.

**Bigger comments, questions, and concerns**

1. I appreciate that you actually implemented your developments in a DA system to see the real-world impacts. However, you only showed results from one forecast, which does not give much confidence regarding the generality or strength of the results. If possible, I strongly urge you to add more cases. I know that adding more cases requires more work, but doing so would not add much to the length of the paper and would make the conclusions much more robust.

We agree that single one case is not convincing to show the capability of our DA system, though the major purpose is to develop the DA system with the cross-correlation process. We have added more cases in the revised manuscript, all nine flight cases from May 15 to June 14, 2010, are applied to DA experiments.

The details of these nine cases are in Table 2 and Figure 8. The results are showed in Figure 11 and Figure 12. The averaged improvement of DA for these nine cases is lower than that for the case on June 3, 2010. The main reasons are that the flight tracks are relative fewer in some cases, or the flight tracks are not around Los Angeles. Another reason is that the initial fields of Control experiment are consistent with observations, especially for the initial time at 00 UTC and 18 UTC. In the case that limited improvements were obtained at the initial fields, the improvements of subsequent forecasts are also low.

2. In my opinion, you neglected to discuss another (and easier) method of dealing with cross-correlations between aerosol species: ensemble-based DA methods (such

as the ensemble Kalman filter) that naturally handle cross-correlations. Thus, I strongly believe you should mention ensemble DA methods in the introduction, and you should cite and briefly discuss Pagowski and Grell (2012) and Schwartz et al. (2014), who assimilated aerosol observations, including PM2.5, with ensemble-based DA methods. There are other references that have also assimilated aerosol observations with ensemble DA, but I believe those two are the most relevant, and without this material regarding ensemble DA, I believe your work is not placed within its proper context.

Thanks. These two papers have been cited, and some discussions about ensemble DA methods are added in the revised manuscript (lines 109-111).

3. In light of the above comment, I believe your title should be more specific, and I suggest adding the word "variational" before "data assimilation".

We agree. The title has been revised.

4. You left out a few important details about the DA system. For example, what DA system were you using? Was it GSI or some other system? Please briefly explain somewhere in the text. Additionally, for your 24-hr forecast you described in section 5, what was the background for DA? Finally, please briefly state the observation errors that you used.

Firstly, this DA system is not GSI or some other widely used systems. It was developed by Li et al. (2013) for the MOSAIC scheme of WRF/Chem model. A simple description about the DA system was added in Section 1 (line 129-130).

Secondly, the backgrounds for DA are the forecasting results from the previous runs without DA. These previous forecasting results have been obtained when we run the model for the BEC statistics. We added the description of the background in Section 5.1 (lines 471-473).

Thirdly, we assume that the observation error is the half of background errors. And a vertical profile of observation errors was applied, resulted from the average of background errors of every level. We think it is an enough large error, even the representativeness error is considered. Since the purpose of this manuscript is to demonstrate the signification of balance constrains in the 3DVAR system, the observation error has an insignificant impact on the analysis of balance constrains. We added the description of the observation error in Section 5.1 (lines 473-475).

5. I believe some aspects regarding Eqs. (6-13) need clarification.
a) Page 8, line 9: Please clarify what you mean by "first variable".
The "first variable" means this variable is fixed. There is not unbalanced component for this variable that is similar to the variable of the vorticity in the DA system of ECMWF (Derber and Bouttier, 1999). But, we did not find the name of "first variable"

in other relevant literatures. Thus, we have removed this name in the revised manuscript.

b) Page 8, Eq. (7): Please fill-in the upper triangle of **K**. Are all upper-triangle elements zero?

Yes, all upper-triangle elements are zero. We have filled in the matrix.

c) Page 8, line 19: Please clarify what you mean by "a one regression coefficient."

It is a mistake. It should be "a regression coefficient". We have revised this sentence.

d) Some more details about how you compute $\rho ij$ would be beneficial.

The $\boldsymbol{\rho_{ij}}$ is the statistical regression coefficients between the variables $i$ and $j$. For example $\boldsymbol{\rho_{12}}$ is the regression coefficient between $\delta EC$ and $\delta OC$. Here, $\delta EC$ and $\delta OC$ are estimated from the forecast differences of 24 h and 48 h forecasts of one month (May 15 to June 14, 2010), that is $\boldsymbol{\delta EC = EC^{24} - EC^{48}}$, $\boldsymbol{\delta OC = OC^{24} - OC^{48}}$. Similar to the calculation of the BEC, $\boldsymbol{\delta EC}$ and $\boldsymbol{\delta OC}$ are also estimated by this forecast difference to represent the difference between real state and forecasts. Using the forecast difference, we have 30 pairs $\delta EC$ and $\delta OC$ for each grid. Then, we can estimate a regression equation and obtain the regression coefficient $\boldsymbol{\rho_{21}}$. Since the $\boldsymbol{\rho_{21}}$ of each grid are close, we use the data of $\boldsymbol{\delta EC}$ and $\boldsymbol{\delta OC}$ at all grids to estimate a regression equation and obtain a regression coefficient $\rho_{21}$. This $\rho_{21}$ should be more robust. Figure 1 shows the scatter plots of $\boldsymbol{\delta EC}$ and $\boldsymbol{\delta OC}$ for all grids. The size of $\boldsymbol{\delta EC}$ or $\boldsymbol{\delta OC}$ is $(N \times 30)$, where N is the number of model grid points, and 30 represents 30 days. From this scatter plots, we can obtain the regression equation:

$$\widehat{\boldsymbol{\delta OC}} = 0.9 \times \boldsymbol{\delta EC}. \tag{1}$$

$\widehat{\boldsymbol{\delta OC}}$ is the predict of $\boldsymbol{\delta OC}$, that is $\boldsymbol{\delta OC_b}$. The residuals are $\boldsymbol{\delta OC_u}$. In this equation, the intercept is neglected, since $\delta OC$ and $\delta EC$ are forecast differences that can be considered to be zero mean values. Using $\boldsymbol{\delta NO_3}$ $\boldsymbol{\delta EC}$ and $\boldsymbol{\delta OC_u}$, we can estimate the regression equation to predict $\boldsymbol{\delta NO_{3b}}$, and obtain $\rho_{31}$ and $\rho_{32}$.

$$\boldsymbol{\delta NO_{3b}} = \widehat{\boldsymbol{\delta NO_3}} = \rho_{31}\boldsymbol{\delta EC} + \rho_{32}\boldsymbol{\delta OC_u}. \tag{2}$$

Then, the other regression equations and regression coefficient can be obtain step by step. Some more detail about the calculation of $\boldsymbol{\rho_{ij}}$ was added in the revised manuscript (lines 327-331).
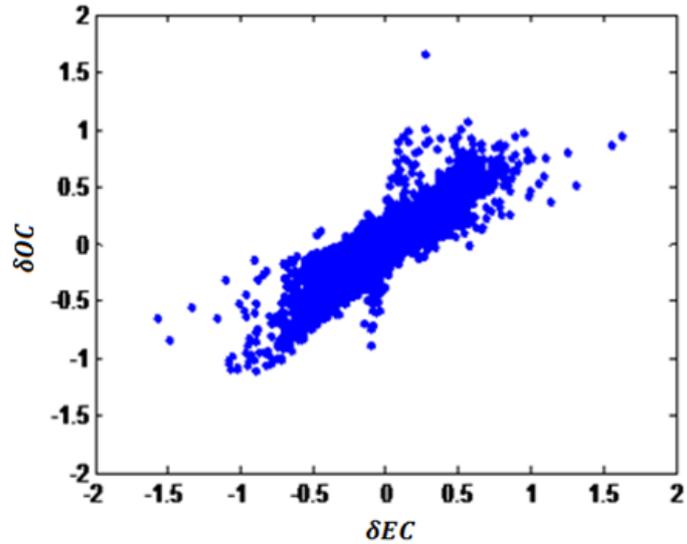
Figure 1 Scatter plots of $\delta EC$ and $\delta OC$

e) Additionally, I think it would be nice if you provided some details on how to interpret $\rho ij$ to bolster the discussion on page 14.

We agree. Some details of the calculation of $\rho_{ij}$ are added in the revised manuscript (lines 327-331).

f) What would happen if the regression was not "based" on EC? In other words, what would happen if you listed the control vector species in reverse [such that OTR was in the first row on the LHS of Eq. (7) and EC was in the last row]? You mention some of this on page 14 lines 6-8, but I believe a clear description about the "order" or "first and second variables" would be greatly beneficial. You also mention using OTR as the "last variable" (page 14, line 19), but the rationale for this choice is not obvious to me. Please clarify.

We think it is difficult to clarify this question which is beyond the scope of current study. We set this order of species mainly due to the following two reasons. First, the correlation of EC and OC is the highest. Second, OTR is correlative with all other variables. The purpose of the balance constrain is to obtain as independent variables as possible. So, EC and OC should be the first two for their high correlation. If we set the other variable such as OTR as the first order, and OC as the second order, the coefficient of determination of the regression equation of OC will be less, compared with the coefficient of determination of the regression equation with EC as the first order. It will increase the correlation of $OC_u$ with the other variables. Similarly, since OTR includes many species that are correlative with former variables, the coefficient of determination of the regression equations of OTR will be largest using all former variables as factors, and then obtain the more independent $OTR_u$. To investigate the impacts by using different orders, more tests need to be conducted which we may address in later studies.

8

In addition, we can refer other DA system to understand this question. In the formulation of DA of ECMWF (Derber and Bouttier, 1999), The balance operator **K** matrix which transforms $[\zeta, \eta_{\mathrm{u}}, (T, P_s)_{\mathrm{u}}, q\,]$ into $[\zeta, \eta, (T, P_s),\ q\,]$. The first variable is the vorticity $\zeta$ The balanced part of the divergence $\eta$ and the temperature and surface pressure $(T, P_s)$ are given by the equations:

$$\eta_{\mathrm{b}} = \boldsymbol{M}\zeta, \tag{3}$$

$$(T, P_s)_{\mathrm{b}} = \boldsymbol{N}\zeta + P\eta_{\mathrm{u}}. \tag{4}$$

Then, the **K** matrix becomes:

$$K = \begin{bmatrix} I & 0 & 0 & 0 \\ M & I & 0 & 0 \\ N & P & I & 0 \\ 0 & 0 & 0 & I \end{bmatrix}. \tag{5}$$

In this DA system of ECMWF, the $\zeta$ and $\eta$ are the first two variables. We think the reason is that they are relative high correlative. The $(T, P_s)$ is the third variable, since it is correlative with the former variables. The $q$ is the last variable that is not correlative with the other variables. Unfortunately, Derber and Bouttier (1999) did not explain why they set this order of variables. We explain it from our thought.

In another reference about the study of balance constraints for GSI system (Chen et al., 2013), the order of control variables is the stream function ($\psi$), the unbalanced part of the velocity potential ($\chi_{\mathrm{u}}$), the unbalanced part of temperature ($T_{\mathrm{u}}$), the unbalanced part of surface pressure ($ps_{\mathrm{u}}$), and the relative humidity ($rh_{\mathrm{u}}$). Here, $rh_{\mathrm{u}}$ is the last order, and its regression equation uses all former variables as factors.

g) Page 9, line 1: I feel like the word "deduced" to describe $\rho21$ is inaccurate. How exactly are you obtaining $\rho ij$?
The "deduced" has been changed to "obtained" (line 362). Please see the response of d) for the process of obtaining $\rho ij$.

h) Page 9, lines 3-10: I found $\boldsymbol{\varepsilon}$ confusing and also unnecessary. By definition, $\boldsymbol{\varepsilon} = \delta OCu$, so why not just use $\delta OCu$ directly in place of $\boldsymbol{\varepsilon}$? Thus, I suggest removing all instances of $\boldsymbol{\varepsilon}$.
We wrote Eq. (8) and reserve residual $\varepsilon$ because it is a normal format for a regression equation. This may be easier to understand for the reader. And we revised the Eq. (9) and its explanation to understand easily the calculation of $\boldsymbol{\delta OC_b}$ and $\varepsilon$ (line 215).

i) Page 9, Eq. (11). I believe you're missing "δ" on EC and OCu.
Corrected, thank you.

**Smaller comments, questions, and concerns**

1. Page 2, line 14: Please clarify what you mean by "coincident".

<span style="color:red">We mean the PM2.5 concentrations of the experiment with balance constraints are more consistent with the observed concentrations. The sentence has been removed.</span>

2. Page 2, line 17: Please omit the word "significant" because you did not perform any statistical significance testing, and you only showed results from one forecast.

<span style="color:red">Thanks, the sentence has been revised.</span>

3. Page 2, line 21: Again, omit "significantly".

<span style="color:red">Thanks, the sentence has been revised.</span>

4. Page 2, line 26: Technically, the observation errors also determine the analysis increments.

<span style="color:red">Thanks, the observation error has been added in the revised manuscript (line 65).</span>

5. Page 3, line 5: Most models now have a state size $O(10^7)$. Suggest modifying.

<span style="color:red">Thanks, the sentence has been revised (lines 69).</span>

6. Page 3, lines 3-8: Note that with ensemble DA methods, these issues are not as difficult to deal with.

<span style="color:red">Thanks, we add the qualifier of "variational data assimilation system" (line 99).</span>

7. Page 3, line 12: Please define in words what you mean by PM2.5.

<span style="color:red">Thanks, the definition of PM2.5 has been added in the revised manuscript (lines 105).</span>

8. Page 3, line 13: Suggest spelling out GSI and adding a reference.

<span style="color:red">We have spelled out GSI in the revised manuscript (line 77).</span>

9. Page 4, lines 9-11: Do these assumptions only apply to variational approaches?

<span style="color:red">Yes. We add the qualifier of "variational" in the revised manuscript (line 99).</span>

10. Page 4, line 20: This might be a good place to mention Pagowski and Grell (2012) and Schwartz et al. (2014).

<span style="color:red">Thanks, we have cited these two references at lines of 109-110.</span>

11. Page 5, line 1: Please spell out "AOD".

<span style="color:red">Thanks, the sentence has been revised (line 119).</span>

12. Page 5, line 4: Please clarify what you mean by "not adjacent".

"that are not ADJACENT" means that are not the connecting. This sentence has been revised (line 121-122).

13. Page 5, line 10: Please clarify what you mean by "eight/four".

The MOSAIC scheme offers flexibility in specifying the number of size bins, four or eight bins are commonly used. Four bins used are located between 0.039–0.1 μm, 0.1–1.0 μm, 1.0–2.5 μm, and 2.5–10 μm. Some introductions of "eight/four" size bins have been added in the revised manuscript (lines 127-128).

14. Page 5, line 12: Suggest "developed" rather than "estimated".

Thanks, the sentence has been revised.

15. Page 6, Eq. (1): It should be J(x) not J(δx).

Thanks, the sentence has been revised (line 159).

16. Page 6, lines 20-25 and Eq. (2): You've ignored non-linear H and its linearization about the background to derive the linear H. In Eq. (1), H is nonlinear, but in Eq. (2) it's linear, because you've linearized H about xb. Please be more precise.

Thanks. In this paper, the observation variables are the species concentration and total PM2.5 concentration. They are really linear relationship with the state variables. Anyway, we have added the assumption of linear in the revised manuscript (line 163).

17. Page 7, line 2: In the expression for the innovation, here H should be nonlinear (H).

Since the relationship between observation variables and state variables is linear. We do not emphasize the nonlinear H.

18. Page 7, line 7: Again, it should probably be $10^7$ rather than $10^6$. Also, $10^{12}$ should probably be $10^{14}$.

Thanks, the sentence has been revised (lines 174).

19. Page 7, line 14: Please clarify what you mean by "is commonly simplified with vertical levels."

The standard deviation matrix (D) is a diagonal matrix with the size of $(N \times m)^2$, that is, each species at each grid has a value of standard deviation. But, to reduce the computational cost, we use the average value of standard deviations that are at the same vertical level. Though the size of **D** is fixed, the number of parameters of standard deviations reduces in the DA system. We have added some introduction in the revised manuscript (lines 180-181).

20. Page 10, lines 8-9: It was unclear to me how you got Eq. (17) from Eq. (6).

Please add some steps or clarify.

According the definition of the BEC,

$$B = \langle(\delta x)(\delta x^T)\rangle.$$

Using Eq. (6),

$$B = \langle(K\delta x_u)(K\delta x_u)^T\rangle$$
$$= \langle(K\delta x_u)(\delta x_u{}^T K^T)\rangle$$
$$= KB_u K^T$$

Some explains have been added in the revised manuscript (lines 237-242).

21. Page 11, line 1: In Eq. (20), it appears you used $\delta x = B^{1/2}\delta z$. Thus, I believe line 1 on page 11 should read $\delta z = B^{-1/2}\delta x$ (note the negative sign on the exponent of B). Please double-check.

Corrected, thank you.

22. Page 12, line 5: Should be "horizontal grid spacing" not "resolution"…they mean different things.

Thanks, the sentence has been revised (lines 274).

23. Page 12, lines 10-12: Please clarify what you mean by "former forecast". Additionally, where do the initial meteorological conditions come from? Are these also from NARR?

The initial meteorological condition is from NARR. For a 48-hour or 24-hour forecast running, we update the initial meteorological condition using the reanalysis NARR data. But for the initial aerosol condition, since there are not reanalysis data, we use the forecast condition from former forecast as the initial condition. The explaination has been added in the revised manuscript (line 281-283).

24. Page 12, line 26 and page 13, line 2: I wonder if you might want to rename "E_ORG" to "E_OC" and "E_PM25" to "E_OTR" to be consistent with the nomenclature of the control variables. If so, please also change on the relevant figure (Fig. 2) and elsewhere in the text.

Since the emission variables are different with the model control variables. The former include many aerosol precursors such as E_SO2, E_NO2. These aerosol precursors can transform into aerosol through chemical process. Thus, the emission variable E_SO4 or E_NO3 are not completely corresponding to the control variables. We use the name of emission variables, consistent with the name in user's guide of WRF/Chem.

25. Page 13, line 5: "With the exception" is misleading and suggests that the diagonal correlations will be < 0.5. Please modify.

We have modified the corresponding sentence in the revised manuscript (Line 305).

26. Page 13, line 9: Suggest "high" rather than "close".
Corrected.

27. Page 14, line 1: I believe it should be Eqs. (6-13) rather than Eqs. (6-12).
Corrected.

28. Page 14, line 2: I believe Eq. (7) is more correct than Eq. (6).
Corrected.

29. Page 14, lines 9-19: Should the control variables here have subscripts "u"? I'm not sure. Please double-check.
We have double-checked.

30. Page 14: Just a comment—I really like Fig. 3.
Thanks.

31. Page 15, line 2: Suggest "obtained" rather than "performed".
Corrected.

32. Page 15, lines 2, 3, and 8: In all these locations, it should be Fig. 4, not Fig. 2.
Corrected, thanks.

33. Page 15, lines 10-11: I believe OTR and NO3 should be OTRu and NO3u, respectively.
Corrected.

34. Page 15, lines 10-11: Please clarify with what the "decreases" are with respect to.
Thanks, we have modified the sentence (Line 369)

35. Page 15, line 17: I believe it should be Eq. (22), not Eq. (21).
Corrected.

36. Page 15, lines 18-25: Please explain how you get the horizontal correlation scale (Ls) from Fig. 5. Is Ls defined as an e-folding distance? Overall, I was a bit confused by your description of Ls —please clarify.
We assume that the decline curve of horizontal correlations is according to the Gaussian function (Fig. 5). Then the intersection of the decline curve and the line of $e^{-\frac{1}{2}} (\approx 0.61)$ can be approximately as the value of horizontal correlation scale. The introduction has been added in the revised manuscript (lines 384-386).

37. Page 15, line 27: I believe OC, NO3, SO4, and OTR should have subscript "u".
Corrected.

38. Page 16, line 1: Please clarify what you mean by "common factors in regression equations".
The common factors mean EC, $OC_u$, and $NO_{3u}$. For example, EC is used four times in the regression Eqs. (6-13), $OC_u$ is used three times, $NO_{3u}$ is used two times. But, it may be puzzling to readers. We have revised this sentence in the manuscript (lines 390-391).

39. Page 16, lines 4-16: Similar to my above comment, please explain how you get the vertical correlation length-scales from Fig. 6.
For the vertical correlation, we use the real values calculated from the forecasting differences in the DA system, but not approximate values from an alternative function. The name of "vertical correlation length-scale" is just a conception to explain the difference between the unbalanced variables and full variables. We have added some explanations in the revised manuscript (lines 397-399, 405-406)

40. Page 16, lines 13-16: I only see very small differences regarding the vertical correlations between the full and unbalanced variables. Perhaps you may wish to modify the text.
The differences of vertical correlation are slight, compared with the difference of horizontal. The main reason is that the vertical correlations are generally affected by the atmospheric boundary layer height. Thus, all vertical correlation decreases rapidly for the level above the boundary layer height. We have added this explanation in the revised manuscript (lines 410-413).

41. Page 17, line 23: Please clarify that DA-balance assimilates the same observations as "DA-full".
For the DA-full experiments and DA-balance experiments, we use the same observation for the data assimilation. This sentence has been revised (lines 469-470).

42. Page 17, line 25: "WRF" not "WRf".
Thanks, the sentence has been revised.

43. Page 17, line 28: Please clarify what you mean by "the initial time".
The initial time means the start time of the model running, which is listed in Table 2.

44. Page 18, lines 4-26: I feel like this discussion slightly misses the main points. In

my opinion, the main point is that the balance constraints can allow observations of a specific species to impact other variables. Even with PM2.5 observations, because the model-simulated PM2.5 is a function of all the control variables, the individual species' fields are adjusted through the BECs, even without a direct observation of the individual species. Thus, without multivariate correlations, an aircraft observation of OC can only impact OC (because the forward operator for OC is only a function of OC), but with the multivariate BECs, an OC observation can now impact OTR or EC. Perhaps you might wish to clarify some aspects of the text along these lines.

Yes. For the BECs without balance constraints, the observation of OC can only impact OC. The crossing effects among species from the BECs with balance constraints. This section has been revised (lines 484-489).

45. Page 19, lines 1-5: I don't believe it is appropriate to describe the smaller RMSEs as "improvements". You're simply looking at fits to observations, which, on their own, do not tell you anything about the relative goodness of your DA system.

The comparison between the analysis PM2.5 against the assimilated observations is known as "sanity check". It can demonstrate the capability of the DA system. In the revised manuscript, we use more data from all nine cases to demonstrate the effects of the DA system. This section has been revised (lines 506-517).

46. Page 19, line 17: The description here of Fig. 11 is incorrect.
Corrected.

47. Page 19, line 19: It should be Fig. 11a, not 1a.
Corrected.

48. Page 19, line 20: Omit "significantly". You can maybe replace it with "substantially".
Thanks, the sentence has been revised.

49. Page 20, lines 16-17: Please clarify what you mean by these lines.
This sentence means the horizontal correlation scales of unbalanced variables are closer than that of full variables. And the vertical correlation scales show similar trend. The sentence has been revised (lines 556-558).

50. Page 20, line 27: Please clarify what you mean by "mutual spread".
The "mutual spread" has been changed to "crossing spread" (line 489).

51. Page 21, line 6: I don't agree with this line. You're only looking at the analysis fits, which does not mean your analysis fields are necessarily better.

This sentence has been removed in the revised manuscript.

52. Page 21, line 20: Please clarify what you mean by a "universal balance constraint".

The balance constraint in this paper is just a statistical relationship. We hope to find a universal balance that can describe the physical or chemical balanced relationship of aerosol variables, similar with the balance constraint of geostrophic balance or temperature-salinity balance in meteorological or oceanic data assimilation. The sentence has been revised in the manuscript (lines 585-586).

53. Table 1: Suggest also pointing to Eq. (7) in the caption. Also, you should annotate the various species on this figure somehow, because it's difficult to look back to Eq. (7).

Thanks, we followed this suggestion.

54. Fig. 2 caption: Suggest "NEI05" rather than just "NEI"

Thanks, the sentence has been revised.

55. Fig. 4 caption: In my opinion, this figure isn't that close to Fig. 3 so I suggest elaborating.

Corrected, thank you .

56. Fig. 5 caption: Suggest pointing to Fig. 4 rather than Fig. 3.

Corrected.

57. Fig. 6: Suggest adding labels of "Height" to the axes.

Corrected.

58. Fig. 7: Suggest adding a unit (meters) to the colorbar.

Corrected.

59. Figs. 8 and 9: The labels above/below the panels are very small. Can these be enlarged?

Corrected.

**Background error covariance with balance constraints for aerosol species and applications in <mark>variational</mark> data assimilation**

Zengliang Zang[1], Zilong Hao[1], Yi Li[1], Xiaobin Pan[1], Wei You[1], Zhijin Li[2] and Dan Chen[3]

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

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

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

Jun 8, 2016

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

**Abstract**

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

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

## 1. Introduction

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

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

One important role that the BEC plays in meteorological data assimilation is to spread information between different variables to produce balanced analysis fields, which employ 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 $PM_{10\text{-}2.5}$ (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 PM2.5. 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 $PM_{2.5}$ 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 $PM_{2.5}$; 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
$$J(x) = \frac{1}{2}(x - x^b)^T \mathbf{B}^{-1}(x - x^b) + \frac{1}{2}(y - \mathbf{H}x)^T \mathbf{R}^{-1}(y - \mathbf{H}x).$$
(1)

160    Here, $x$ is the vector of the state variables, including EC, OC, $NO_3$, $SO_4$ and OTR; $x^b$ is the

161    background vector of these five species, which are generated by the MOSAIC scheme; $y$ 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 $y$;

164    and $\mathbf{B}$ is the background error covariance associated with $x^b$. Eq. (1) is usually written in the

165    incremental form

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

167    where $\delta x$ ($\delta x = x - x^b$) is the incremental state variable. The observation innovation vector is

168    known as $d = y - \mathbf{H}x^b$. The minimization solution is the analysis increment $\delta x$, and the final

169    analysis is $x^a = x^b + \delta x$. 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), $x^b$ 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 $x^b$ 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 $$\mathbf{B} = \mathbf{DCD^T},\tag{3}$$

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 <mark>To reduce the computational cost, we use the average value of standard deviations that are at the</mark>

181 <mark>same level.</mark> Thus, the standard deviation is simplified with vertical levels. $\mathbf{C}$ is a symmetric

182 matrix, having the form

183 $$\mathbf{C} = \begin{bmatrix} \mathbf{C_{EC}} & \mathbf{C_{EC}^{OC}} & \mathbf{C_{EC}^{NO_3}} & \mathbf{C_{EC}^{SO_4}} & \mathbf{C_{EC}^{OTR}} \\ \mathbf{C_{OC}^{EC}} & \mathbf{C_{OC}} & \mathbf{C_{OC}^{NO_3}} & \mathbf{C_{OC}^{SO_4}} & \mathbf{C_{OC}^{OTR}} \\ \mathbf{C_{NO_3}^{EC}} & \mathbf{C_{NO_3}^{OC}} & \mathbf{C_{NO_3}} & \mathbf{C_{NO_3}^{SO_4}} & \mathbf{C_{NO_3}^{OTR}} \\ \mathbf{C_{SO_4}^{EC}} & \mathbf{C_{SO_4}^{OC}} & \mathbf{C_{SO_4}^{NO_3}} & \mathbf{C_{SO_4}} & \mathbf{C_{SO_4}^{OTR}} \\ \mathbf{C_{OTR}^{EC}} & \mathbf{C_{OTR}^{OC}} & \mathbf{C_{OTR}^{NO_3}} & \mathbf{C_{OTR}^{SO_4}} & \mathbf{C_{OTR}} \end{bmatrix},\tag{4}$$

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. <mark>They represent the correlation</mark>

186 <mark>among pairs of grid points for one species</mark>. 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 $\mathbf{C_{OC}^{EC}} = (\mathbf{C_{EC}^{OC}})^{\mathbf{T}}$. 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 <mark>between different species</mark> 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 ($\delta x_b$) and unbalanced components ($\delta x_u$):

194 $$\delta x = \delta x_b + \delta x_u.\tag{5}$$

195     Note the EC does not need to be divided. <mark>There is not unbalanced component for EC</mark> 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 ($\delta x_u$) to full variables ($\delta x$) by the balance operator $\mathbf{K}$

199 is given by

200 $$\delta x = \mathbf{K}\delta x_u .\tag{6}$$

201  Eq. (6) can be written as

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

203  where $\boldsymbol{\rho}_{ij}$ 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 $\boldsymbol{\rho}_{ij}$ 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 $\boldsymbol{\rho}_{ij}$ is a constant value for all grid points, which are denoted as $\rho_{ij}$

207  and are calculated by linear regression with all grid points. For example, $\rho_{21}$ can be obtained

208  from the regression equation of OC and EC as

209
$$
\boldsymbol{\delta OC} = \rho_{21}\boldsymbol{\delta EC} + \boldsymbol{\varepsilon},
\tag{8}
$$

210  where $\boldsymbol{\varepsilon}$ is the residual. $\boldsymbol{\delta EC}$ and $\boldsymbol{\delta OC}$ 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 $\boldsymbol{\delta EC}$ and $\boldsymbol{\delta OC}$ represent forecast differences that

213  can be considered to be zero mean values. After obtaining $\rho_{21}$, the balanced part (e.g., the value

214  of the regression prediction) of $\boldsymbol{\delta OC}$ can be obtained by

215
$$
\boldsymbol{\delta OC_b} = \widehat{\boldsymbol{\delta OC}} = \rho_{21}\boldsymbol{\delta EC}.
\tag{9}
$$

216  Where $\widehat{\boldsymbol{\delta OC}}$ represents the predicted value of Eq. (8), which is equal to the balanced part ($\boldsymbol{\delta OC_b}$).

217  Remove the $\boldsymbol{\delta OC_b}$ from the full variables to obtain the unbalanced part ($\boldsymbol{\delta OC_u}$), that is, $\boldsymbol{\varepsilon}$ in Eq.

218  (8). Thus, the calculation of $\boldsymbol{\delta OC_u}$ can be written as

219
$$
\boldsymbol{\delta OC_u} = \boldsymbol{\delta OC} - \rho_{21}\boldsymbol{\delta EC}.
\tag{10}
$$

220  Here, $\boldsymbol{\delta OC_u}$ and $\boldsymbol{\delta EC}$ are employed as predictors in the next regression equation to obtain

221  $\boldsymbol{\delta NO_{3b}}$. Then, we can obtain the unbalanced parts of the remaining variables, which are defined as

222  follows:

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

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

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

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

227     can be expressed as

228
$$R^2 = \frac{\text{SSR}}{\text{SST}}, \tag{14}$$

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
$$\mathbf{B_u} = \mathbf{D_u C_u D_u^T}, \tag{15}$$

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
$$\mathbf{C_u} = \begin{bmatrix} \mathbf{C_{EC}} & & & & \\ & \mathbf{C_{OCu}} & & & \\ & & \mathbf{C_{NO_{3u}}} & & \\ & & & \mathbf{C_{SO_{4u}}} & \\ & & & & \mathbf{C_{OTRu}} \end{bmatrix}. \tag{16}$$

237 According the definition of the BEC,

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

239 And $\mathbf{B_u}$ can be written as

240
$$\mathbf{B_u} = \langle (\boldsymbol{\delta x_u})(\boldsymbol{\delta x_u}^T) \rangle. \tag{18}$$

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

242
$$\mathbf{B} = \mathbf{K B_u K^T}. \tag{19}$$

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
$$\mathbf{B}^{\frac{1}{2}} = \mathbf{K B_u^{\frac{1}{2}}}. \tag{20}$$

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

247
$$\mathbf{B}^{\frac{1}{2}} = \mathbf{K D_u C_u^{\frac{1}{2}}}. \tag{21}$$

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

249     state variable:

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

250

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

$$J(\delta z) = \frac{1}{2}\delta z^T \delta z + \frac{1}{2}\left(\mathbf{H}\mathbf{K}\mathbf{D_u}\mathbf{C_u}^{\frac{1}{2}} \delta z - d\right)^T \mathbf{R}^{-1}\left(\mathbf{H}\mathbf{K}\mathbf{D_u}\mathbf{C_u}^{\frac{1}{2}} \delta z - d\right). \tag{23}$$

252

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

$$\mathbf{C_u} = \mathbf{C_{ux}} \otimes \mathbf{C_{uy}} \otimes \mathbf{C_{uz}}. \tag{24}$$

255

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

$$\mathbf{C_u}^{\frac{1}{2}} = \mathbf{C_{ux}}^{\frac{1}{2}} \otimes \mathbf{C_{uy}}^{\frac{1}{2}} \otimes \mathbf{C_{uz}}^{\frac{1}{2}} \tag{25}$$

262

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 unspeciated 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 $\delta EC$，$\delta OC$，

328    $\delta NO_3$ $\delta SO_4$ and $\delta OTR$. 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 $\rho_{ij}$

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 $\rho_{ij}$ in Eq. (7))

| species | regression coefficient ($\rho$) | | | | | 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.



(a) full variables                    (b) unbalanced variables

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, $NO_3$, $SO_4$ 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

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



(a) full variables    (b) unbalanced variables

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

For the correlation matrix of $\mathbf{C}$ and $\mathbf{C_u}$, they are factorized as three independent one-dimensional correlation matrices in Eq. (24). The horizontal correlation $\mathbf{C_x}$ or $\mathbf{C_y}$ is approximately expressed by a Gaussian function. The correlation between two points $r_1$ and $r_2$ 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 for $\mathbf{C_x}$ and $\mathbf{C_y}$, which are considered to be isotropic (Li et al., 2013). This scale can be estimated 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.



(a) full variables        (b) unbalanced variables

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 PM2.5 monitoring sites existed in the innermost domain of the

428 WRF/Chem model (Fig. 7). <mark>The second type of observation data is the speciated concentration</mark>

429 <mark>along the aircraft flight track. The aircraft observations were investigated from the California</mark>

430 <mark>Research at the Nexus of Air Quality and Climate Change (CalNex) field campaign in 2010. Nine</mark>

431 <mark>flights data around Los Angeles from 15 May to14 June, 2010 are selected as the cases of data</mark>

432 <mark>assimilation. Table 2 shows the start time and end time of each flight.</mark> The species of the aircraft

433 observations include OC, $NO_3$, $SO_4$ and $NH_4$. Note that $NH_4$ is not a control variable; thus, the

434 aircraft observation of $NH_4$ 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 <mark>Table 2 The periods of flight during CalNex 2010 and the initial time of assimilation</mark>

| 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 |

441

(a) 00:00 UTC±1.5 h, May 17

(b) 18:00 UTC±1.5 h, May 19

(c) 18:00 UTC±1.5 h, May 21

(d) 00:00 UTC±1.5 h, May 25

(e) 06:00 UTC±1.5 h, May 30

(f) 06:00 UTC±1.5 h, May 31

(g) 12:00 UTC±1.5 h, Jun 02        (h) 12:00 UTC±1.5 h, Jun 03



(i) 18:00 UTC±1.5 h, Jun 14

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 ±1.5h, 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 aircarft 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, $NO_3$, $SO_4$ 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 thee 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, $NO_3$ and

489   $SO_4$ 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.
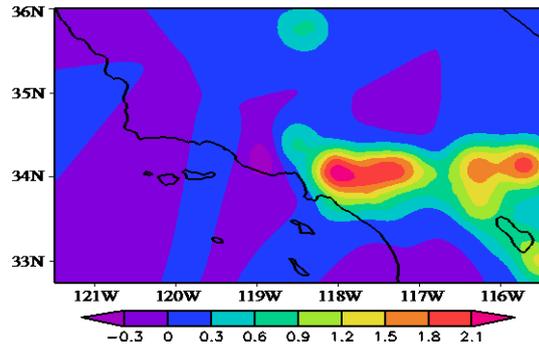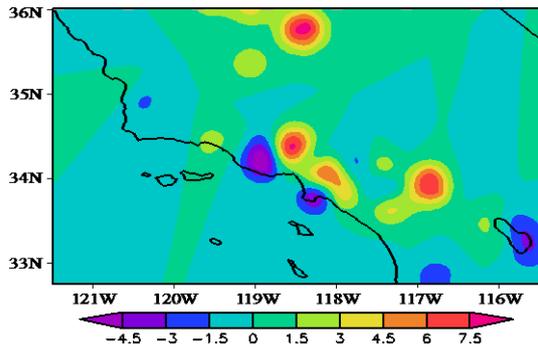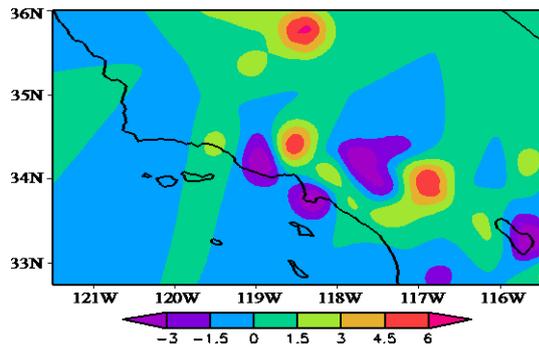
(a) EC in the DA-full

(b) EC in the DA-balance
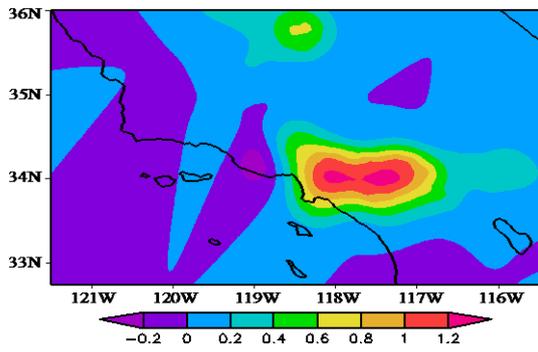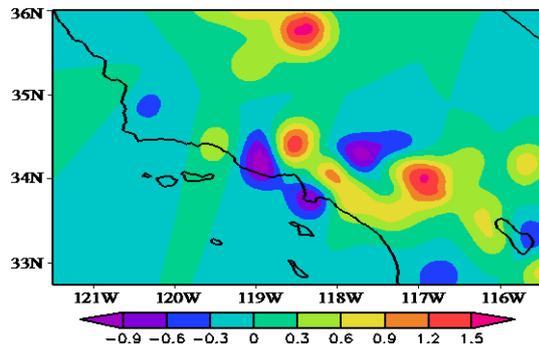
(c) OC in the DA-full
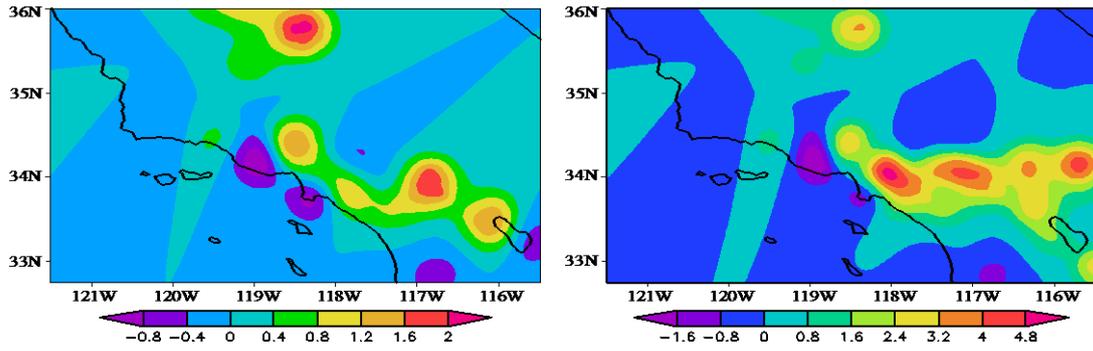
(d) OC in the DA-balance

(e) NO$_3$ in the DA-full

(f) NO$_3$ in the DA-balance
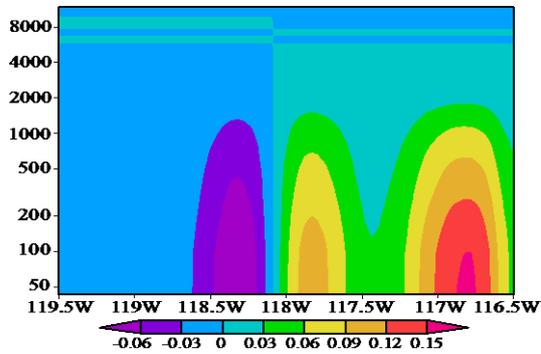
(g) SO$_4$ in the DA-full
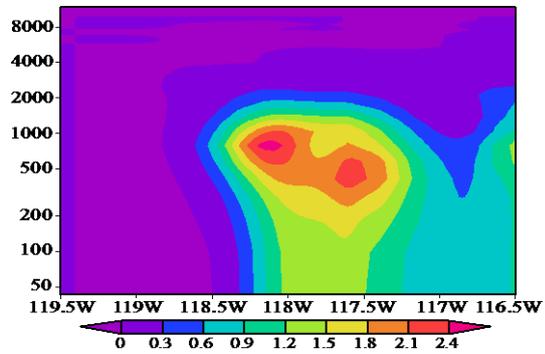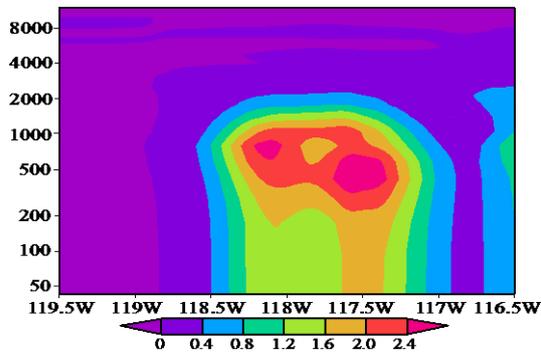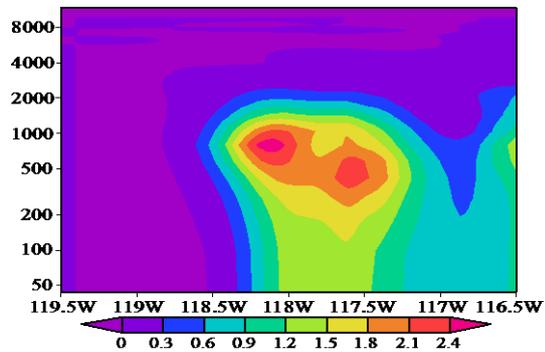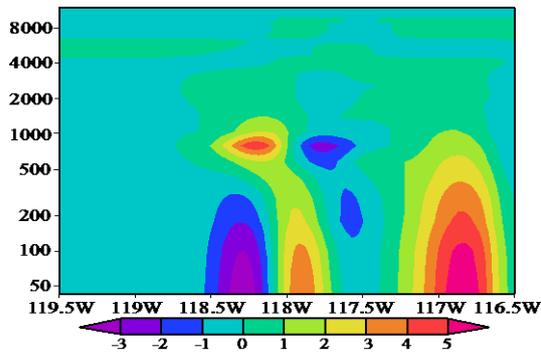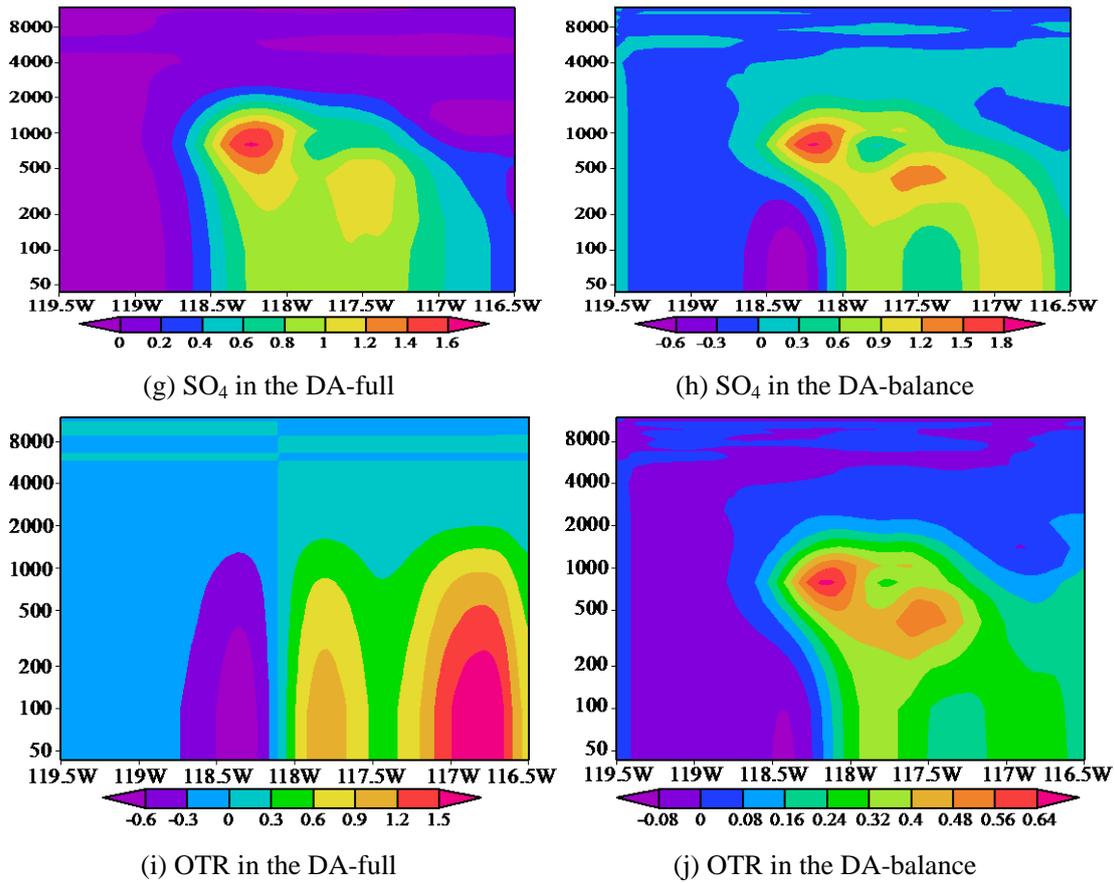
(h) SO$_4$ in the DA-balance

(i) OTR in the DA-full

(j) OTR in the DA-balance

Figure 9. Surface distributions of increments of the five variables of EC, OC, NO$_3$, SO$_4$ and OTR at 12:00 UTC on June 3, 2010. The left column and right column are from DA-full and DA-balance, respectively.



(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) NO$_3$ in the DA-full

(f) NO$_3$ in the DA-balance

(g) SO$_4$ in the DA-full                    (h) SO$_4$ in the DA-balance



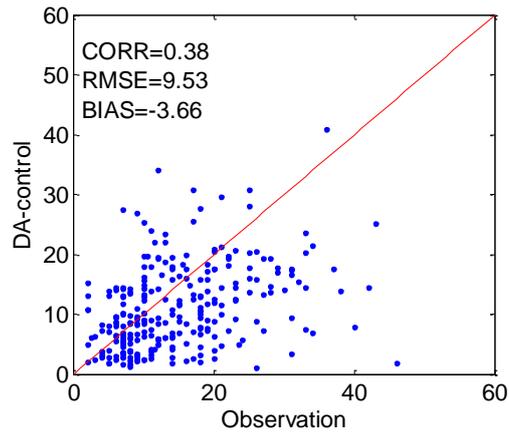(i) OTR in the DA-full                    (j) OTR in the DA-balance

503    Figure 10. Same as Figure 9, with the exception of the vertical sections along 35 N.
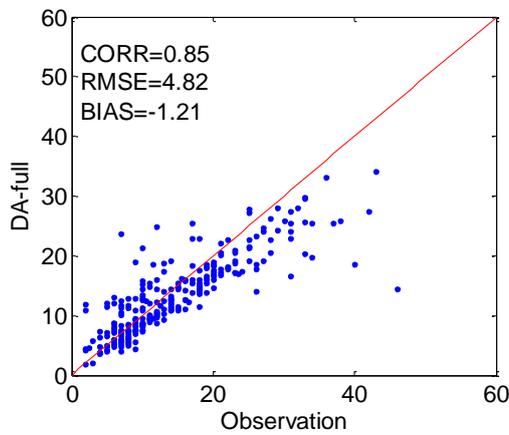
504
505    **5.3 Evaluation of data assimilation and forecasts**

506    <mark>Figure 11 shows the scatter plots of the initial model fields versus the surface observation for all</mark>

507    <mark>nine cases. In Fig. 11a, the simulated concentrations of the Control experiment display a</mark>

508    <mark>significant underestimation with a BIAS of -3.66μg/m$^3$. The mean concentration of Control is</mark>

509    <mark>10.90 μg/m$^3$, about 25.1% lower than observed mean concentrations (14.56 μg/m$^3$). In the DA-full</mark>

510    <mark>and DA-balance experiments, there are remarkable increases for the simulated concentrations, and</mark>

511    <mark>the BIASs reduce to as small as -1.21 and -0.94 μg/m$^3$. The RMSE is 9.53 μg/m$^3$ in the Control</mark>

512    <mark>experiment. The RMSE reduces to 4.82 and 4.48 μg/m$^3$ in the DA-full and DA-balance</mark>

513    <mark>experiment, respectively. There are also significant improvements for the CORR in the DA-full</mark>

514    <mark>and DA-balance experiments, compared with the Control experiment. Furthermore, these three</mark>

515    <mark>statistical measures of the DA-balance experiments show some slight improvement, compared</mark>

516    <mark>with that of the DA-full experiments. The result demonstrates that more observation information</mark>

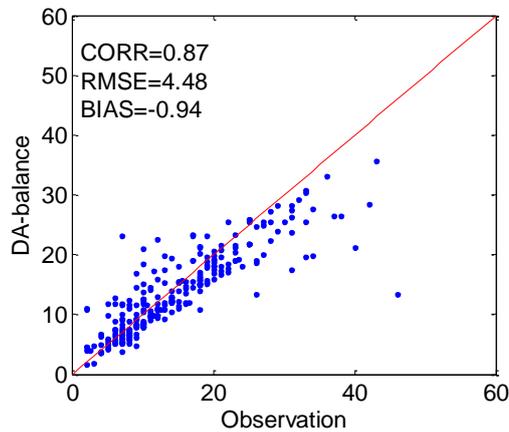517    <mark>spread by balance constraints can improve assimilation performance.</mark>

(a) Control



(b) DA-full
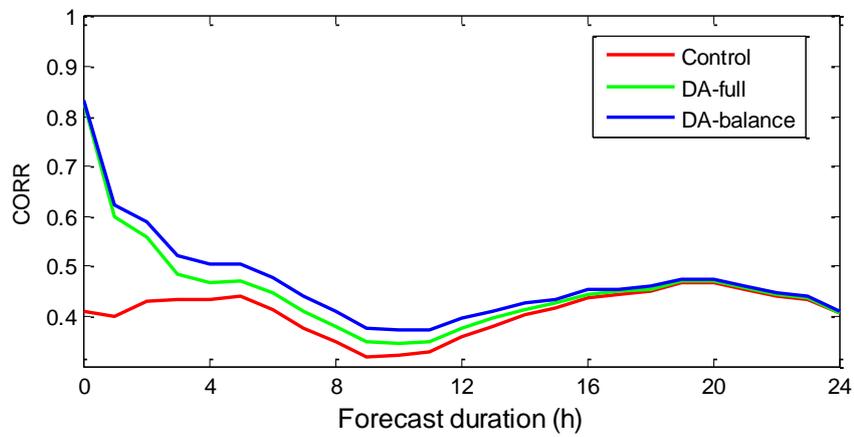


(c) DA-balance

Figure 11. Scatter plots of observed concentrations of $PM_{2.5}$ versus simulated $PM_{2.5}$ concentrations of the experiments of (a) Control, (b) DA-full, and (c) DA-balance for all nine cases.
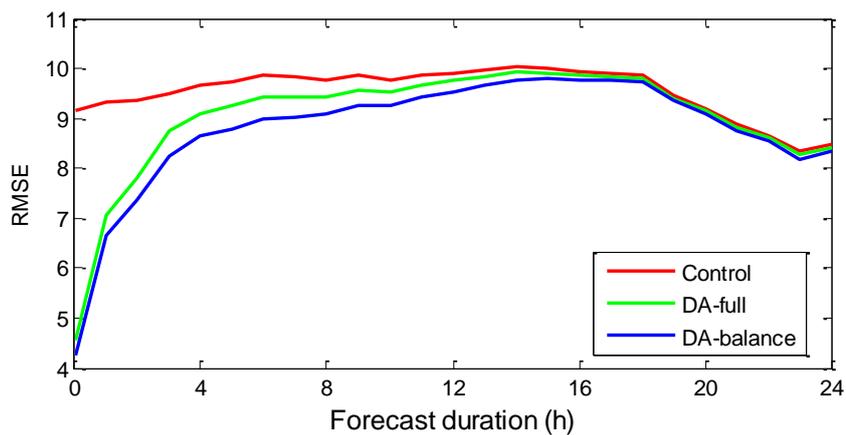
To evaluate the effects of the data assimilation, the CORR, RMSE and BIAS during the forecast time are calculated for each case, and their averaged results are showed in Figure 12. The CORRs of the DA-balance and DA-full experiments are very close (Fig. 12a). But, the difference increase
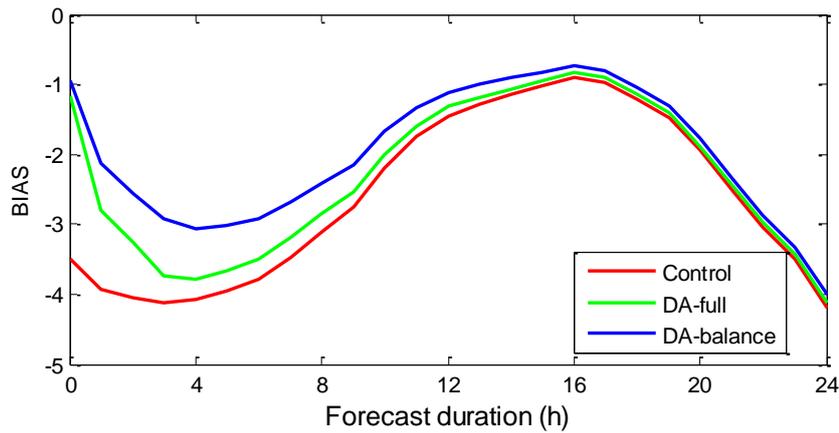
after the first hour with a higher CORR in the DA-balance experiment. The CORR of the DA-balance experiment is substantially higher than that of the DA-full experiment from the 2[nd] hour to the 16[th] hour. Similar improvements for the RMSE and the BIAS of the DA-balance experiment are observed in Fig. 12b and Fig. 12c, respectively. The improvement for the BIAS in the DA-balance experiment is the most significant among these three statistical measures. The peak value of the improvement for the BIAS (Fig 12c) is at the 4[th] hour, and the improvement is distinct until the end of forecasts. These improvements indicate that the balance constraint is positive for the subsequent forecasts, which derives from the balanced initial distribution among species.



(a) CORR



(b) RMSE

(c)  BIAS

**6. Summary and discussion**

We examined the BEC in a 3DVAR system, which uses five control variables (EC, OC, NO$_3$, SO$_4$ and OTR) that are derived from the MOSAIC aerosol scheme in the WRF/Chem model. Based on the NMC method, differences within a month-long period between 24- and 48-h forecasts that are valid at the same time were employed in the estimation and analyses of the BEC. The background errors of these five control variables are highly correlated. Especially between EC and OC, their correlation is as large as 0.9.

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

To evaluate the impact of the balance constraints on the analyses and forecasts, three groups of 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, $NO_3$ and $SO_4$ 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 $2^{nd}$ hour to the $16^{th}$ 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 ($NO_2$), sulfur dioxide ($SO_2$),

588    and (carbon monoxide) CO. These gaseous pollutants are correlated with some aerosol species,

589    such as $NO_3$, $SO_4$ 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 availablity**

594    This data assimilation system is established by ourself. The code of this system can be obtained on

595    request from the first author (zzlqxxy@163.com).

596

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 PM2.5 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. Colrado: NOAA Earth System Research Laboratory.

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

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

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

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

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

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

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

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

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

Schwartz, C. S., Liu, Z., Lin, H.-C., Cetola, J. D., 2014, Assimilating aerosol observations with a "hybrid" variational-ensemble data assimilation system, J. Geophys. Res. Atmos., 119, 4043–4069, 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

Table 1 Regression coefficients of balance operator K and the coefficient of determination

(regression coefficients correspond to $\rho_{ij}$ in Eq. (7))

| species | regression coefficient ($\rho$) | | | | | 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

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, $NO_3$, $SO_4$ 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 $e^{-\frac{1}{2}} (\approx 0.61)$ 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, $NO_3$, $SO_4$ 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 μg/m$^3$) and (c)

760    mean bias (BIAS in μg/m$^3$) of the $PM_{2.5}$ concentration forecasts against observations as a function

761    of forecast duration.

762