

# Ensemble-Trained $PM_{2.5}$ Source Apportionment Approach for Health Studies

DONGHO LEE,\*,†,‡  
 SIVARAMAN BALACHANDRAN,‡  
 JORGE PACHON,‡ ROSHINI SHANKARAN,‡  
 SANGIL LEE,§ JAMES A. MULHOLLAND,‡  
 AND ARMISTEAD G. RUSSELL‡

*Gyeongnam Province Institute of Health and Environment, Changwon, Gyeongnam 641-702, Korea, School of Civil and Environmental Engineering, Georgia Institute of Technology, Atlanta, Georgia 30332-0512, and Center for Analytical Measurement Services, Korea Research Institute of Standards and Science, Daejeon 305-340, Korea*

*Received February 13, 2009. Revised manuscript received June 30, 2009. Accepted July 15, 2009.*

An ensemble-trained chemical mass balance (CMB) approach is developed for particulate matter (PM) source apportionment (SA), particularly for use in health studies. The approach uses results from a short-term emission-based chemical transport model (CTM) and multiple receptor-based approaches. Ensemble results have less day-to-day variation in source impacts and fewer biases between observed and estimated  $PM_{2.5}$  mass compared to the original receptor model results. Ensemble results show increases in road dust, biomass burning, and coal impacts, but secondary organic carbon (SOC) impacts decrease. These results, along with observations, are then used to obtain new source profiles. Two sets of new source profiles based on ensemble results in summer (July 2001) and winter (January 2002) were developed, and used in separate CMB applications for a 12-month data set of daily  $PM_{2.5}$  measurements at the Atlanta, GA, Jefferson Street site. Results show that ensemble-trained CMB approaches, using both summer profiles and winter profiles, effectively reduce day-to-day variability of source impact estimates by reducing fewer days of zero impact from sources known to be present, as compared to traditional receptor modeling, suggesting improved results.

## 1. Introduction

There are a number of reasons to develop direct relationships between emission sources, air pollutant concentrations, and health end points. In addition to establishing associations between sources and atmospheric pollutant concentrations (1–4), uncertainties in epidemiologic analyses might be reduced using particular matter (PM) source impacts rather than PM components (2, 5), and modeling approaches that link physically or chemically transformed products to sources may decrease classification error in epidemiologic analysis. However, there are issues associated with use of current source apportionment (SA) approaches in health analyses. As explored in Marmur et al. (2), the more commonly used

receptor models, as well as those involving physically based models, have limitations that affect their usefulness. For example, commonly applied receptor models do not include all sources. Baek et al. (6) found that a typical receptor model application includes/identifies sources/factors representing only 60–80% of inventoried primary emissions. Incomplete specification of the sources leads to both errors and bias. Receptor model SA approaches appear to introduce excessive day-to-day variability, as suggested by major sources with zero impact on some days and significant variation in relative impacts from one day to the next. A further limitation is that the information is representative of the location of the measurements. There are further issues specific to the method employed as well. For example, the chemical mass balance (CMB) (7) approaches depend on accurate characterization of source emissions and typically do not account for species transformations. Factor analytic approaches, such as Positive Matrix Factorization (PMF) (8, 9), on the other hand, result in factors which can include contributions from multiple sources and can be difficult to interpret. In addition, calculated impacts can change from subjective choices in the model application. Receptor modeling approaches are also limited in their ability to deal with species transformations, and hence in their ability to tie PM health impacts to sources of secondary pollutants.

These limitations of receptor modeling approaches suggested using detailed, physically based, chemical transport modeling (CTM), such as the Community Multiscale Air Quality (CMAQ) model (10). Marmur et al. applied CMAQ to the southeastern United States (2) with the intent of employing directly calculated source impacts on air quality in time-series health studies. Such an application is computationally demanding. Another identified limitation was that the fraction of total PM contributed by each source varied very little day by day. The lack of variation in source fractions was attributed to the fact that (1) the meteorological model used did not fully capture the fine scales (time and space) variations (11), (2) emission inputs had little day-to-day variation, except for weekends and changes due to meteorology, and (3) emissions are allocated based on land-use, not capturing actual spatial–temporal variations. On the other hand, the analysis suggested that the CMAQ results might be appropriate for use for longer term exposures (e.g., to study chronic health effects), and have the advantage of large, continuous spatial coverage. A critical issue for both receptor-based and emission-based SA approaches is how to formally evaluate the results. Emission-based approaches can be indirectly evaluated by assessing how well the model correctly simulates the temporal and spatial concentrations of the observed air pollutants. Good agreement between model predictions and measurements, however, does not guarantee that the predicted source impacts are correct as there may be compensatory errors. Evaluating receptor model results is even more problematic. Models can be tested on artificial data sets, but that raises the question of how realistic such sets are in terms of capturing all of the complexities that impact PM composition at a receptor (e.g., source profile variation, source strength variation, physical and chemical transformations, and unidentified source contributions) (12). Alternatively, models can be tested on “event” days in which impacts from individual sources are known to dominate. Such events, however, are unusual, and can cause problems when applying factor analytic methods. Comparison of results from various modeling approaches can provide a false sense of comfort because the models may share the same assumptions or mathematical limitations. The more apparent

\* Corresponding author phone: 404-894-7693; 404-894-8266; e-mail: estlake@gmail.com.

† Gyeongnam Province Institute of Health and Environment.

‡ Georgia Institute of Technology.

§ Korea Research Institute of Standards and Science.

hybridizations of the methods are also flawed, e.g., inverse modeling (13, 14) or using a chemical transport modeling as a smart interpolator. The former still lacks day-to-day variability; the latter retains the increased, likely excessive, variability.

One approach that can help overcome the limitations discussed above is to use ensembles of model results (15, 16). In this case, results from various models and approaches, along with model uncertainties, can be used together. This work explores such an approach, and extends it one step to develop source profiles for later use.

## 2. Approach

This work relies on blending receptor-based modeling and detailed and short-term emission-based CTM results to then use to “train” a CMB application.

**Chemical Mass Balance (CMB) Approaches Including CMB-RG, -MM, -LGO.** The CMB method is widely used to estimate contributions of emission sources to receptors. Weighted least-squares solutions are solved to an equation, which describes a mass balance between ambient concentration of each chemical species and a linear sum of products of source profiles and source contributions:

$$C_i = \sum_{j=1}^m f_{ij} \cdot S_j + e_i \quad (1)$$

where  $C_i$  is the observed concentration of species  $i$ ,  $f_{ij}$  is the fraction of species  $i$  in source  $j$ ,  $S_j$  is the contribution of source  $j$ , and  $e_i$  is for error. The fractional compositions of species in each source (source profiles) and ambient concentrations of each species at the receptor serve as input data to the CMB model. Typically, a weighted least-squares approach is used to find the source contribution ( $S_j$ ) that minimizes the differences between calculated and observed species concentrations. Uncertainty in input data is used both to weigh the contribution of input data and to calculate the uncertainty of the output data. Outputs are the impacts of the modeled sources at the receptor and their uncertainties. Historically, observed metals, ions, elemental (EC), and organic carbon (OC) PM concentrations are used. This type of application, in this paper referred to as CMB-RG (regular), is typically conducted using the effective variance-based approach (17) available from U.S. EPA (e.g., CMB-8 (18)).

Another CMB approach uses organic molecular markers, referred to here as CMB-MM, and allows separate identification of more sources of primary organic PM than are typically elucidated using CMB-RG. While providing such additional source impacts, the need for additional chemical analysis limits its application. This approach has been applied by Cass and co-workers in Atlanta and elsewhere (19–21). Again, the effective variance approach is used.

CMB typically uses fixed source profiles which are based on measurements that may have significant errors when applied at different specific locations and times from which the original tests were performed. Source emissions change with time, so the measured profiles may be out of date. In addition, there can be substantial variation in source profiles by region, season, and ambient conditions (e.g., changes in gasoline composition and engine operating characteristics between summer and winter may result in different profiles). Marmur et al. (22) extended CMB to allow a constrained variation of source profiles, including constraints based on source profile uncertainties as well as gaseous species measurements. This approach, called CMB-LGO (Lipschitz Generalized Optimization), reduced but did not eliminate some obvious errors in traditional CMB results, including random peaks and zeroes in source impacts. In their CMB-

LGO study, Marmur et al. used measurement-based source profiles (MBSP) from several different studies; source profiles used for gasoline vehicle (GV) and diesel vehicle (DV) emissions were from NFAQS (23); profiles for biomass burning (BURN) and coal combustion (COAL) were from BRAVO (24); profiles for dust (DUST) were from more regionally representative measurements in Alabama (25).

**Factor Analytic (FA) Approaches.** A second class of receptor model is based on FA principles, with the most commonly applied approaches being principle component analysis (PCA), PMF (9, 26), UNMIX (27), and multilinear engines (MEs) (23, 24). PMF and MEs have seen increasing application, particularly by Paatero, Hopke and co-workers (including to the Atlanta Jefferson Street (JST) SEARCH and STN sites), and a new version is now available from U.S. EPA (28).

For FA and CMB approaches, incomplete specification of all the sources contributing to PM concentrations at a receptor can cause bias and increased variability (2, 15, 26); that is, with a limited number of factors identified or source profiles available, these methods assign mass from other sources to available factors/sources. Further, source impacts covary as they are strongly impacted by meteorology. Increased variability arises from source/factor profiles being collinear (to varying degrees), so using SA results can increase classification error in epidemiologic modeling.

**Emissions-Based Approaches Using a Chemical Transport Model (CTM).** Emissions-based source apportionment modeling relies on the use of a CTM, such as CMAQ (10, 29, 30), and requires emissions (e.g., SMOKE (31)) and meteorological inputs (e.g., MM5 (32)). Here we use results from CMAQ extended to include a direct source apportionment technique by adding tracers (CMAQ-TR) (6) that follow 33 sources of primary PM species. This has been applied to a domain covering the contiguous United States, with a primary focus on the Southeast, for comparison with detailed observations from the JST site.

**Model Ensemble. Compilation and Comparison of PM Source Apportionment Results.** Source impacts at JST were calculated by combining CMB-RG, CMB-MM, CMB-LGO, PMF, and CMAQ-TR results from prior work (e.g., CMAQ-TR (6), CMB-RG (33), PMF (33), CMB-MM (19)), or those calculated here (CMB-LGO). The study used data from the JST site in Atlanta because it contains daily speciated PM<sub>2.5</sub> data, including metals, ionic species, EC/OC, and detailed organic speciation for the periods studied here. Details of the measurements at JST are available elsewhere (34, 35). Furthermore, source apportionments have been conducted using JST data in prior studies, though the method is designed to be applicable at any site where multiple source apportionment methods can be applied (which includes any routine monitor with speciated PM data). However, the main goal of this study was to test how well this approach works and we chose this site for testing. Future studies will include applications at other sites. Results from CMAQ-TR were available for July 2001 and January 2002 and based on inputs developed as part of the 2002 VISTAS modeling, which included day-specific emissions from major point sources. The Regional Planning Organization (RPO) consensus domain, which was used here, provides coverage over the U.S. and parts of Canada and Mexico (Figure S1). CMAQ-TR (6) provides temporal and spatial source impacts of 33 primary source categories, and source contributions to ionic species.

In this study, we apportion PM<sub>2.5</sub> with nine pollutant sources: GV, DV, DUST, COAL, BURN, secondary organic carbon (SOC), SULFATE, NITRATE, and AMMONIUM. This includes about 80% of the inventoried primary PM<sub>2.5</sub> sources and the dominant secondary sources. “Vehicle source” impacts used in CMB-RG and PMF methods are divided into GV and DV, based on their relative impact ratio in the three

other SA approaches; the three nonprimary sources, ammonium sulfate, ammonium bisulfate, and ammonium nitrate, used in CMB-RG, CMB-LGO, and PMF, are divided into SULFATE, NITRATE, and AMMONIUM sources based on molecular weights. Out of 33 primary sources (Table S1) used in CMAQ, 17 were recategorized to five primary sources that were also used in CMB and PMF applications: (1) BURN impacts include agricultural burning, cigarettes, fireplaces, prescribed burning, yard waste burning, wood combustion, and wild fires; (2) DUST includes construction, and paved and unpaved road dust; (3) GV includes gasoline-fuel engine sources such as gasoline vehicles, nonroad aircraft, and nonroad gasoline engines; (4) DV includes diesel vehicles and nonroad diesel engines; and (5) COAL includes just the primary PM from coal burning. However, the dominant precursor of sulfate PM in the eastern U.S. is the emission of SO<sub>2</sub> from the combustion of coal.

*Ensemble-Derived Source Impacts:* Daily PM source impact results using the different SA approaches at the JST site in Atlanta were ensembled using weighted averaging:

$$\bar{S}_j(t_k) = \frac{\sum_{i=1}^L w_{ji}(t_k) \cdot S_{ij}(t_k)}{\sum_{i=1}^L w_{ji}(t_k)} \quad (2)$$

where  $\bar{S}_j(t_k)$  is the ensemble-calculated impact of source  $j$  (in ug/m<sup>3</sup>) at time  $t_k$ ,  $S_{ij}(t_k)$  is the impact developed by method  $i$  (e.g., CMAQ, CMB-RG, PMF, etc.), and  $w_{ji}(t_k)$  is the weight assigned to the calculated source impact as described below. Ensembling was applied for two different periods: July 2001 for summer and January 2002 for winter. There were 2–5 different  $S_{ij}(t_k)$  available for each day during July 2001 and January 2002. Uncertainties in each  $S_{ij}(t_k)$  were assessed and the inverse values were used for the weights.

$$w_{ji} = \frac{1}{\sigma_{S_{ij}}} \quad (3)$$

We explored using both  $(1)/(\sigma_{S_{ij}})$  and  $(1)/(\sigma_{S_{ij}}^2)$  and chose  $(1)/(\sigma_{S_{ij}})$  because we found a lower average reduced chi-square value (hereafter referred to as chi-square value) and fewer zero-impact days using  $(1)/(\sigma_{S_{ij}})$ , particularly for DUST. In future work, we will explore the impact of using different weights and ways of estimating uncertainties to optimize performance based on average chi-square, zero-impact days, and day-to-day variability.

In CMB, uncertainties in the source contribution are estimated using an effective variance approach and are an output of the model (17):

$$\sigma_{S_j} = \left( \sum_{i=1}^n \frac{f_{ij}^2}{\sigma_{ei}^2 + \sum_{j=1}^m \sigma_{fij}^2 S_j^2} \right)^{-1/2} \quad (4)$$

where  $\sigma_{S_j}$  (in ug/m<sup>3</sup>) is the uncertainty in source contribution  $S_j$  (in ug/m<sup>3</sup>),  $\sigma_{ei}$  is the uncertainty in the measurement of ambient species  $i$ , and  $\sigma_{fij}$  is the uncertainty in the fraction of  $i$  species, in the  $j$  source profile. Uncertainties calculated using eq 4, available from prior studies, were used for CMB-RG (33) and CMB-MM (19), or calculated here for CMB-LGO and PMF. PMF factor profile uncertainties were developed using bootstrapping (28, 33). For CMAQ, uncertainties were adjusted based on the relations between CMAQ and CMB-RG models (derivation is in Supporting Information):

$$4\sigma_{S_j}^2 = 4\sigma_{S_j,CMB}^2 + (S_{j,CMAQ} - S_{j,CMB})^2 \quad (5)$$

where  $\sigma_{S_j,CMB}$  (in ug/m<sup>3</sup>) is an uncertainty estimate of source impact  $S_j$  (in ug/m<sup>3</sup>) from CMB-RG models, and  $S_{j,CMAQ}$  and  $S_{j,CMB}$  are contribution estimates of source  $j$  in CMAQ and CMB-RG, respectively. This approach directly accounts for differences between observed and simulated concentrations.

*Ensemble-Trained CMB Source Apportionment.* While we view both CMB and EA approaches as attractive for application to long-term PM source apportionment, here we concentrate on using a CMB-based approach, similar to CMB-LGO. One reason to choose CMB-LGO is that it allows for constrained optimization of the source profiles to attain a best fit to ensemble results and can readily utilize gas-phase pollutant concentration data. In this case, CMB-LGO was applied to develop constrained ensemble-based source profiles (EBSPs) that best replicate the ensemble-derived  $\bar{S}_{jk}$ . The general steps in the training are as follows: (1) use a best-fit to the ensemble SA series and develop optimized source profiles (i.e., EBSPs), (2) calculate the new source impacts and assess variability. The EBSPs should decrease the day-to-day variability found in the initial CMB application, and account for location specific differences in relative source profile composition.

First, EBSPs which best replicate the ensemble results are calculated using the species mass balance equation, taking the observed concentrations,  $C_{jk}$ , and ensemble source impacts,  $\bar{S}_{jk}$ , as known, leaving the source profiles,  $f_{ij}$ , as unknown:

$$C_{ik} = \sum_{j=1}^m \bar{f}_{ij} \cdot \bar{S}_{jk} + e_{ik} \quad (6)$$

Equation 6 was solved for two periods using observed data and ensemble results at JST in July 2001 for summer and January 2002 for winter, by minimizing  $\chi^2$ :

$$\chi^2 = \sum_{i=1}^n \frac{(C_{ik} - \sum_{j=1}^m \bar{f}_{ij} \cdot \bar{S}_{jk})^2}{\sigma_{C_{ik}}^2} \quad (7)$$

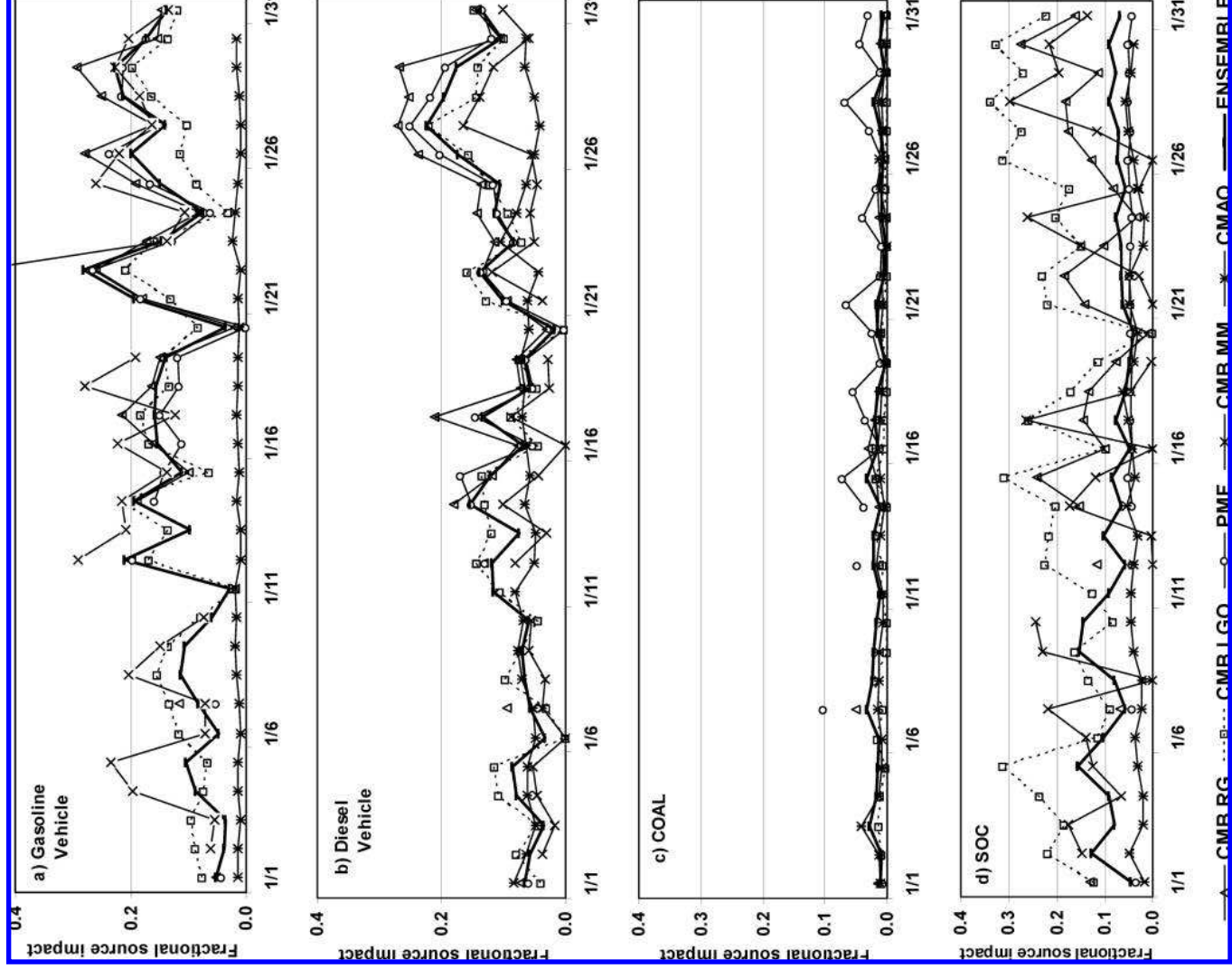
where  $\bar{f}_{ij}$  are the ensemble-based source profiles. For all sources, the sum of relative source profile compositions,  $\sum_{i=1}^N \bar{f}_{ij}$ , is limited to be less than 1 to consider the oxidized forms of metals and the relation between OC and OM (OM/OC  $\geq$  1.4). Lower and upper bounds for the fractions of species, which were developed by Marmor et al. (22), are applied for the profile optimization, except the OC/EC ratio bounds for GV and DV sources, which are adopted from Cadle et al. (36), and the upper bounds for SO<sub>4</sub><sup>2-</sup> and OC species in COAL source, which are adopted from the BRAVO study (Table S2) (24). Source impacts for each day in 2002, using the EBSPs in CMB-LGO,  $S_{jk}^*$ , were then recalculated taking the  $C_{ik}$  and the  $\bar{f}_{ij}$  as known. Gas phase pollutant concentrations of CO, SO<sub>2</sub>, and NO<sub>y</sub> provide constraints to the solutions as discussed in Marmor et al. (35).

There is no reason that the ensemble method cannot be easily employed at other sites, and it is not critical to have results from any one specific method. If measurement data are available, it is straightforward to apply CMB-RG, CMB-LGO, and PMF. Since the CMAQ domain includes all of the contiguous U.S., most receptor sites will be covered by a single CMAQ run as long as model runs are available for the corresponding time of the ensemble (e.g., Jan 2002 and July 2001, as done here, and many groups have applied CMAQ to annual and even multiyear periods). However, CMB-MM is the most resource-intensive method and therefore is not expected to be as widely utilized.

## Results

**Ensemble.** Daily PM source apportionment results, developed by five different SA methods and the ensemble-





**FIGURE 1. Comparison of daily fractional source impact, developed from the ensemble approach, to the SA results from traditional models: (a) GV, (b) DV, (c) COAL, (d) SOC.**

derived SA, show that the ensemble technique reduces the day-to-day variability of source impacts to varying degrees and decreases the number of zero-impact days for sources such as diesel exhaust and coal combustion (Figure 1). Daily variability is reduced when compared to CMB-LGO results for GV, COAL, and SOC sources, while little change is found for other sources. There were no zero-impact days with the ensemble averaging, while the original CMB-LGO results showed three zero-impact days in July 2001 (two days for COAL, one day for DV) and 11 zero-impact days in January 2002 (nine days for COAL, one day each for DV and SOC).

The number of source impacts available for each day varied for each source since receptor-based SA results were not conducted for those days missing the necessary ambient species measurements. However, in the original CMB-LGO,

missing values were replaced by the geometric mean of the measured values, and their corresponding uncertainties were set equal to four times the geometric mean (35). All sources had all five source impacts on most days (Table S3) in summer and winter except for DUST in winter, which had CMB-MM results for only five days. It should be noted that for both COAL and BURN in the summer and for COAL in the winter, only four SA results were available since CMB-MM did not resolve those sources in these cases.

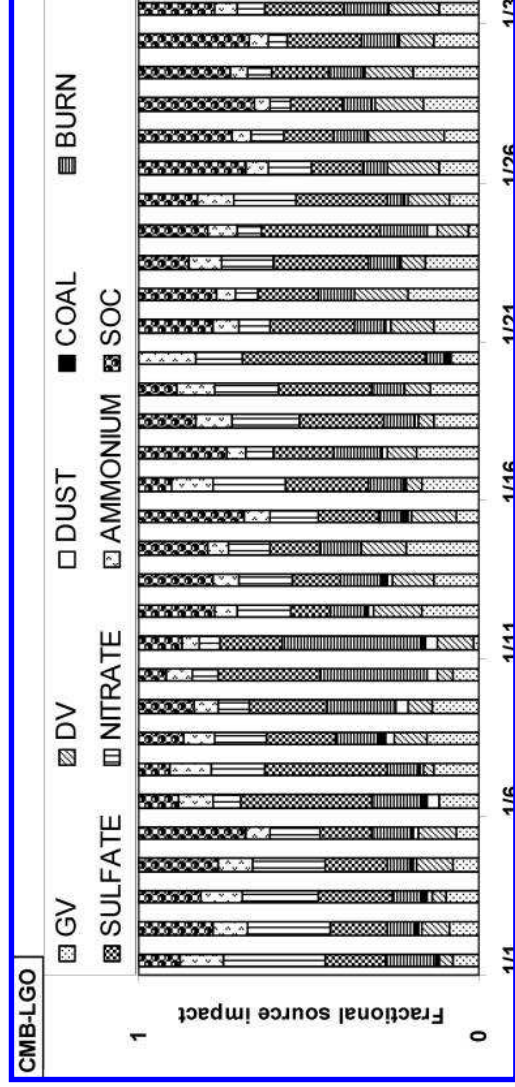
Source impacts of several components increased compared to the original CMB-LGO results. The January fractional source impacts for DUST, BURN, and COAL sources increased to 3.3%, 20.1%, and 1.5% from 1.4%, 12.0%, and 0.5%, respectively. SOC decreased from 19.6% in the original CMB-LGO to 8.1% in the ensemble approach during January 2002 (Table 1). Similarly, the July fractional source impacts for

**TABLE 1. Comparison of Monthly Fractional Source Impact in Ensemble Results with CMB-LGO Results at the Atlanta Jefferson Street (JST) Site for July 2001 and January 2002**

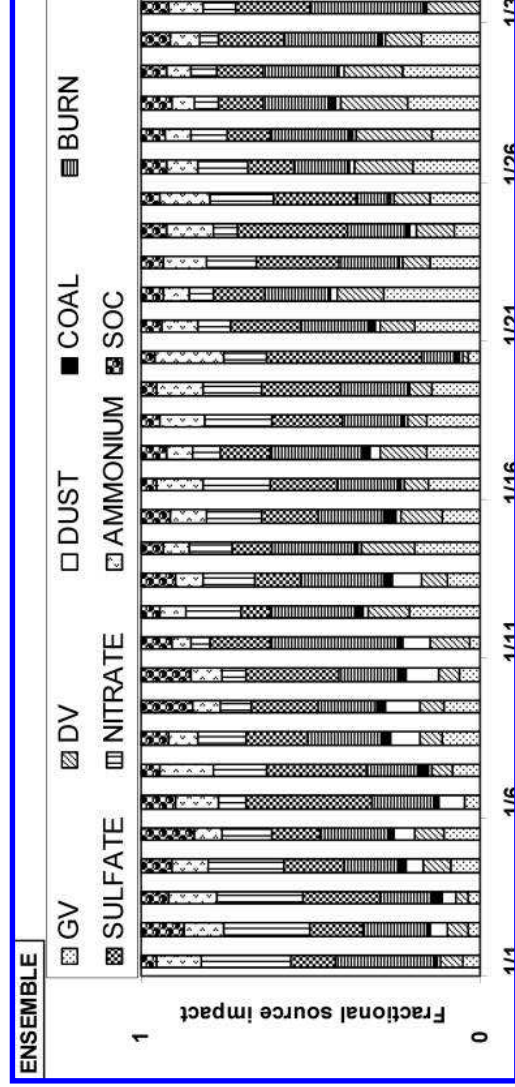
	January 2002		July 2001	
	CMB-LGO	ensemble	CMB-LGO	ensemble
GV	0.121 (0.049)	0.127 (0.067)	0.044 (0.024)	0.035 (0.022)
DV	0.096 (0.050)	0.098 (0.050)	0.060 (0.030)	0.076 (0.026)
DUST	0.014 (0.10)	0.033 (0.032)	0.038 (0.061)	0.054 (0.071)
BURN	0.120 (0.073)	0.201 (0.058)	0.070 (0.056)	0.062 (0.030)
COAL	0.005 (0.006)	0.015 (0.008)	0.005 (0.005)	0.016 (0.015)
SO <sub>4</sub> <sup>2-</sup>	0.226 (0.090)	0.199 (0.080)	0.457 (0.081)	0.458 (0.076)
NO <sub>3</sub> <sup>-</sup>	0.137 (0.063)	0.139 (0.062)	0.030 (0.012)	0.032 (0.014)
NH <sub>4</sub> <sup>+</sup>	0.085 (0.028)	0.106 (0.034)	0.171 (0.025)	0.184 (0.030)
SOC	0.196 (0.082)	0.081 (0.031)	0.124 (0.065)	0.082 (0.026)
$\chi^2$	14.6 (14.0)	11.1 (7.9)	23.4 (22.9)	21.3 (20.7)

DUST and COAL increased, whereas SOC decreased in ensemble results. Increases in primary sources were driven by CMAQ and PMF, which show greater impacts from primary sources than the other SA approaches. The decrease in SOC impacts in the ensemble approach is due to the negative bias in PMF-simulated SOC, which we find estimates less SOC than CMB. In addition to the impact changes between sources, the day-to-day fractional impacts of each source are more stable in the ensemble results (Figures 2 and 3),

particularly the COAL and BURN fractions. The unusually high fraction of DUST for 10 days (2nd–6th, 8th–11th, and 13th of January 2002) (Figure 3) is caused by the absence of CMB-RG, CMB-MM, and PMF results during those days due to missing observations. Therefore, the SA results from only two methods, CMB-LGO and CMAQ, are used for ensemble weighting and, as a result, the higher dust impact in CMAQ dominates. The chi-square value averages, which use the errors between observed and estimated PM<sub>2.5</sub> mass for those two months, are decreased to 11.1 and 21.3 in January ensemble and July ensemble, respectively, from 14.6 and 23.4 in CMB-LGO. One concern appears to be successfully addressed in this approach as day-to-day variability in the receptor-based approach is reduced and modeling performance is improved. We also recalculated the ensemble trained source impacts by removing CMAQ and, separately, CMB-MM to see how sensitive the ensemble procedure was to these two SA methods. We found that overall, the ensemble source apportionments changed by less than 3% (contribution of any one source to the total) when these methods were removed (Figure S2). The major reason that removing CMAQ and CMB-MM had such little impact is that the ensemble averages are driven by CMB-RG, CMB-LGO and PMF, which typically have smaller uncertainties than CMAQ or CMB-MM. This characteristic of ensembling with weights based on uncertainty has a net effect of reducing variability.



**FIGURE 2. Daily fractional contributions to PM<sub>2.5</sub> in SA results developed from CMB-LGO method at JST during January 2002 (35).**



**FIGURE 3. Daily fractional contributions to PM<sub>2.5</sub> in ensemble results at JST during January 2002.**

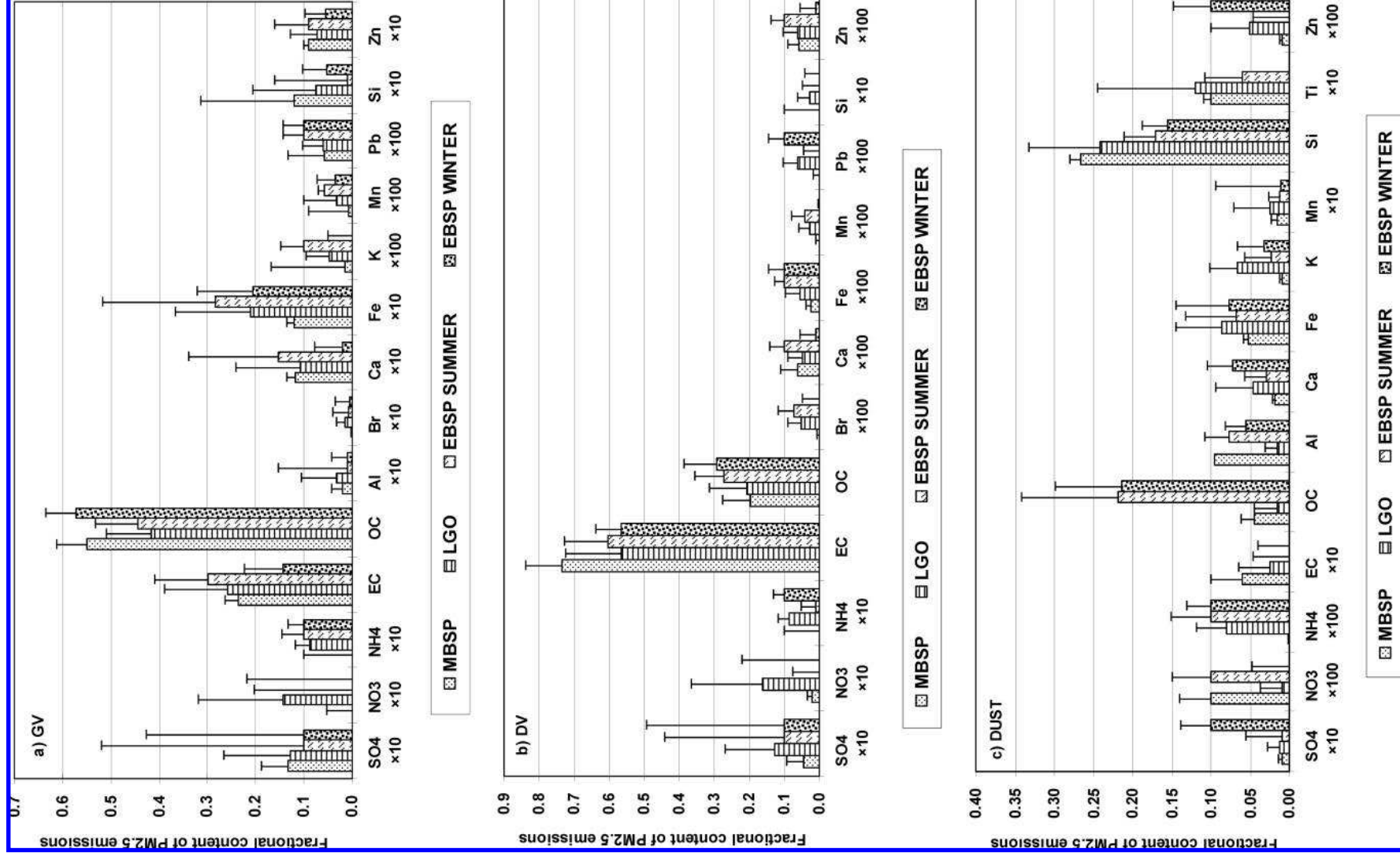


FIGURE 4. Measured (MBSP) (23–25), LGO-derived (22), and ensemble-based (EBSP) source profiles: (a) GV, (b) DV, (c) DUST. Bars represent one standard deviation. The LGO profile was derived from a 447-day data set; the EBSP SUMMER and EBSP WINTER were derived from a 24-day data set in July 2001 and a 21-day data set in January 2002, respectively. (Note: BURN and COAL are shown in Figure S3).



**TABLE 2. Performance Measures of Source Apportionments Using EBSPs and MBSP**

used source profile	summer		winter	
	EBSP	MBSP	EBSP	MBSP
chi-square	7.2 ± 6.9	19.5 ± 9.5	13.2 ± 16.8	26.6 ± 19.6
<sup>a</sup> PM mass ratio, %	88.8 ± 9.7	87.2 ± 9.6	98.3 ± 15.4	93.6 ± 15.1
<sup>b</sup> correlation of species, R	0.994 ± 0.013	0.994 ± 0.013	0.999 ± 0.006	0.999 ± 0.008
frequency of zero source impact, day	8	8	16	25
gasoline vehicle	0.046 ± 0.026	0.060 ± 0.028	0.095 ± 0.047	0.098 ± 0.040
diesel vehicle	0.083 ± 0.041	0.080 ± 0.037	0.101 ± 0.056	0.099 ± 0.045
road dust	0.037 ± 0.038	0.027 ± 0.031	0.014 ± 0.009	0.010 ± 0.009
biomass burning	0.089 ± 0.053	0.075 ± 0.040	0.182 ± 0.085	0.090 ± 0.072
coal burning	0.008 ± 0.007	0.007 ± 0.007	0.015 ± 0.010	0.008 ± 0.008
SULFATE	0.359 ± 0.081	0.363 ± 0.082	0.213 ± 0.085	0.228 ± 0.088
NITRATE	0.057 ± 0.039	0.058 ± 0.040	0.123 ± 0.057	0.129 ± 0.059
AMMONIUM	0.146 ± 0.028	0.147 ± 0.029	0.107 ± 0.032	0.113 ± 0.032
SOC	0.174 ± 0.068	0.182 ± 0.068	0.150 ± 0.067	0.225 ± 0.082

<sup>a</sup> Estimated-to-observed PM mass ratio. <sup>b</sup> Correlation of species between observed and estimated amount.

**Optimized Source Profile Compositions.** Two sets of source profile compositions were calculated using the two different seasonally based ensemble results. Only 21 days in January 2002 and 24 days in July 2001 had all the ambient data (all ions, carbon fractions, metals, and gases) desired, and are used for the profile optimization procedure (Tables S4 and S5). Source profiles of the five primary sources, GV, DV, DUST, BURN, and COAL are calculated, using CMB LGO, to most closely reproduce the ensemble SA results.

Ensemble-based source profiles (EBSP) are similar to those found applying CMB-LGO, though some profile species vary markedly. The GV winter EBSP has an OC/EC ratio of 4.0, which is higher than both the LGO-generated profile of 1.6 and that from NFRAQS of 2.3. However, the OC/EC ratio of 1.5 in the summer EBSP is similar to the other two. The total carbon (TC) fractions in PM<sub>2.5</sub> from gasoline vehicle emission were 0.71 in the winter and 0.74 in the summer EBSP, similar to both 0.78 in NFRAQS and 0.67 in the LGO generated profiles. The Zn fraction (Zn is considered a good marker for gasoline vehicle impact) is 0.005 for winter and 0.009 for summer compared to those from the two other procedures of 0.009 (NFRAQS) and 0.008 (LGO).

The OC/EC ratios in the DV EBSPs, 0.52 in the winter and 0.45 in the summer, differ from the NFRAQS ratio of 0.27, but are within a typical range of 0.3–0.9, as reviewed by Hodan, et al. (37). Their findings suggested that the OC fraction of 0.19 and EC fraction of 0.75 in the current HDDV profile, 1996 EPA PM<sub>2.5</sub> Split Factor Profile Number 35600, might overestimate the EC fraction and underestimate the OC fraction due to a failure to consider variable burn conditions, truck size, fuel type, and control scenarios. TC contents in the DV profiles are stable by season, 0.86 for the winter and 0.88 for the summer EBSP, being similar to 0.93 in the NFRAQS profiles.

Biomass burning source profiles are characterized by high OC/EC ratios and high K fractions. The OC/EC ratio of 3 in both summer and winter BURN EBSPs is lower than the measured profiles of 4.1 from BRAVO and 4.7 from the LGO profiles. Further, the K fractions, which drive the source contribution calculation of biomass burning, are also lower in both the winter EBSP (0.024) and summer EBSP (0.035), than the MBSP (0.056) and LGO (0.063) profiles. These differences may be due to burning conditions, wood type, grate feature, burning rate, etc. For instance, the estimated uncertainties in the vegetative burning profile from the BRAVO study (MBSP in Figure 4) are up to 97% for the fraction of EC and 98% for the K fraction due to the variability in burned materials. On the other hand, the OC and TC fractions in biomass burning

EBSPs do not show much seasonal variation, with the same OC fractions of 0.55 in both summer and winter profiles and TC fractions of 0.88 in the summer and 0.86 in the winter. These values are similar to the OC fraction of 0.64 and TC fraction of 0.80 in the BRAVO profiles.

Fractions of major species in the DUST EBSPs are similar across seasons: 0.21 and 0.22 for OC, 0.08 and 0.07 for Fe, 0.03 and 0.02 for K, 0.16 and 0.17 for Si in the winter and summer EBSPs, respectively. However, the Ca fraction differed by season, being 0.07 in winter and 0.03 in the summer. OC and Fe fractions in the DUST EBSPs are slightly higher than the typical ranges of 0.04–0.19 for OC and 0.02–0.05 for Fe that have been reported in other studies (24, 37). On the other hand, both the MBSPs and the LGO source profiles are in the lower end of this typical range for OC. The higher OC levels in the DUST EBSPs are driven by PMF, which we have found to have OC fractional contents on the order of 0.07–0.45. Our results place the OC level for dust at 0.20, which is slightly higher than the typical range, but still lower than PMF results, and is likely due to road dust, which is higher in OC than soil dust. This also explains the higher Fe levels.

In the COAL EBSPs, sulfate fractions are 0.47 in the winter and 0.13 in the summer. The fractions of OC, EC, Ca, Fe, and K are all lower in winter than in summer in varied levels, while the fraction of selenium is 0.007 in both seasons. The seasonal variation of sulfate appears excessive if we consider that the chemical composition of PM<sub>2.5</sub> emissions from coal burning is theoretically stable by season. This bias comes from using PMF results that show a larger impact from coal burning, which can produce additional sulfate by oxidation of the concurrent SO<sub>2</sub> emissions.

Species fractions in the EBSPs are optimized from the relations between fixed source impacts (ensembled source contribution) and PM<sub>2.5</sub> compositions at the receptor. Thus, theoretically, the EBSPs are expected to best replicate the ensemble source impact, but there is no direct method to estimate the accuracy of the optimized profiles. Instead, performance measures from applying CMB-LGO source apportionment using EBSPs are compared to the original CMB-LGO using MBSP (Table 2).

**Source Apportionment Using EBSPs.** The two seasonal EBSPs are applied to estimate source impacts of each day in 2002 (over 320 days for which all species of ambient data are available) at JST. Source apportionments for the winter season (January, February, November, and December) use winter EBSPs while source apportionments for the other eight months (from March to October) use summer EBSPs. Source impacts developed from ensemble-trained CMB (using

EBSPs) show a decreased number of zero-impact days in major sources, being reduced from 33 days in the original CMB-LGO application to 24 days when using the ensemble-trained CMB (mainly from COAL). Several performance measures (Marmur et al. (22)), assessing the quality of fit achieved in SA modeling suggest improved performance: a lowered chi-square value of 13.2, a higher predicted-to-observed PM<sub>2.5</sub> ratio of 98.3%, and high correlation (R) between observed and estimated species concentrations of 0.999 were achieved in the ensemble-trained CMB approach in winter, compared to 26.6, 93.6%, and 0.999, respectively, in the original CMB-LGO.

Applying EBSPs for PM source apportionment introduces relative contribution changes between sources. During the winter, source impacts for biomass burning, road dust, and coal burning in this approach increased to 18.2%, 1.4%, and 1.5%, respectively, from 9.0%, 1.0%, and 0.8% in the original CMB-LGO approach, while the SOC impact decreased to 15.0% from 22.5%. In summer, the relative source impact changes are not as great as in the winter: the GV impact decreased to 4.6% from 6.0%, but the DUST fraction increased to 3.7% from 2.7% in the LGO results.

Source impact changes replicate changes between the original and ensemble results. However, the ensemble-trained PM<sub>2.5</sub> source apportionment still has issues with regard to errors from unused unavailable primary source profiles such as meat cooking, natural gas burning, and “industrial–mineral,” which are known to contribute to OC in the Southeast (6).

The advantage of the recalculating approach using EBSPs instead of straight ensemble results is that the former can be used for long-term source apportionment. Ensembling for longer periods is limited by a large computational cost in the current emission-based SA method (CMAQ). In addition, as mentioned before, ensemble results can be dominated by a single SA approach for days in which some methods fail to develop SA results, and lead to biased source impact estimates for those days. The recalculating approach using EBSPs addresses the risk of biased source impacts.

## Acknowledgments

This study was supported by the U.S. EPA under Grants RD-83215910, RD83362601, RD83096001, RD82897602, and RD83107601, Southern Company, and Gyeongsangnam-do Province Government of Korea.

## Supporting Information Available

This information is available free of charge via the Internet at <http://pubs.acs.org>.

## Literature Cited

- Laden, F.; Neas, L.; Dockery, D. W.; Schwartz, J. Association of fine particulate matter from different sources with daily mortality in six US cities. *Environ. Health Perspect.* **2000**, *108*, 941–947.
- Marmur, A.; Park, S. K.; Mulholland, J. A.; Tolbert, P. E.; Russell, A. G. Source apportionment of PM<sub>2.5</sub> in the southeastern United States using receptor and emissions-based models: Conceptual differences and implications for time-series health studies. *Atmos. Environ.* **2006**, *40*, 2533–2551.
- Stolzel, M.; Laden, F.; Dockery, D. W.; Schwartz, J.; Kim, E.; Hopke, P. K.; Neas, L. M. Source apportionment of fine and coarse particulate matter and daily mortality in two US cities - A comparison of different methods. *Epidemiology* **2005**, *16* (5), S95–S95.
- Thurston, G. D.; Ito, K.; Mar, T.; Christensen, W. F.; Eatough, D. J.; Henry, R. C.; Kim, E.; Laden, F.; Lall, R.; Larson, T. V.; Liu, H.; Neas, L.; Pinto, J.; Stolzel, M.; Suh, H.; Hopke, P. K. Workgroup report: Workshop on source apportionment of particulate matter health effects - Intercomparison of results and implications. *Environ. Health Perspect.* **2005**, *113* (12), 1768–1774.
- Sarnat, J. A.; Marmur, A.; Klein, M.; Kim, E.; Russell, A. G.; Sarnat, S. E.; Mulholland, J. A.; Hopke, P. K.; Tolbert, P. E. Fine particle

sources and cardiorespiratory morbidity: An application of chemical mass balance and factor analytical source-apportionment methods. *Environ. Health Perspect.* **2008**, *116* (4), 459–466.

- Baek, J.; Park, S. K.; Hu, Y.; Russell, A. G. Source Apportionment of Fine Organic Aerosol Using CMAQ Tracers. In *The 4th Annual CMAS Models-3 Users' Conference, Durham, NC, September 27, 2005*; CMAS: Durham, NC, 2005.
- Watson, J. G.; Zhu, T.; Chow, J. C.; Engelbrecht, J.; Fujita, E. M.; Wilson, W. E. Receptor modeling application framework for particle source apportionment. *Chemosphere* **2002**, *49* (9), 1093–1136.
- Hopke, P. K.; Ramadan, Z.; Paatero, P.; Norris, G. A.; Landis, M. S.; Williams, R. W.; Lewis, C. W. Receptor modeling of ambient and personal exposure samples: 1998 Baltimore particulate matter epidemiology-exposure study. *Atmos. Environ.* **2003**, *37* (32), 4595–4595.
- Paatero, P.; Hopke, P. K.; Hoppenstock, J.; Eberly, S. I. Advanced factor analysis of spatial distributions of PM<sub>2.5</sub> in the eastern United States. *Environ. Sci. Technol.* **2003**, *37* (11), 2460–2476.
- Byun, D.; Schere, K. L. Review of the governing equations, computational algorithms, and other components of the Models-3 Community Multiscale Air Quality (CMAQ) modeling system. *Appl. Mech. Rev.* **2006**, *59* (1–6), 51–77.
- Hogrefe, C.; Porter, P. S.; Gego, E.; Gilliland, A.; Gilliam, R.; Swall, J.; Irwin, J.; Rao, S. T. Temporal features in observed and simulated meteorology and air quality over the eastern United States. *Atmos. Environ.* **2006**, *40* (26), 5041–5055.
- Brinkman, G.; Vance, G.; Hannigan, M. P.; Millford, J. B. Use of synthetic data to evaluate positive matrix factorization as a source apportionment tool for PM<sub>2.5</sub> exposure data. *Environ. Sci. Technol.* **2006**, *40* (6), 1892–1901.
- Gilliland, A. B.; Appel, K. W.; Pinder, R. W.; Dennis, R. L. Seasonal NH<sub>3</sub> emissions for the continental United States: Inverse model estimation and evaluation. *Atmos. Environ.* **2006**, *40* (26), 4986–4998.
- Mendoza-Dominguez, A.; Russell, A. G. Iterative inverse modeling and direct sensitivity analysis of a photochemical air duality model. *Environ. Sci. Technol.* **2000**, *34* (23), 4974–4981.
- Guo, Z. C.; Dirmeyer, P. A.; Gao, X.; Zhao, M. Improving the quality of simulated soil moisture with a multi-model ensemble approach. *Q. J. Royal Meteorol. Soc.* **2007**, *133* (624), 731–747.
- Pereira, M. B.; Berre, L. The use of an ensemble approach to study the background error covariances in a global NWP model. *Mon. Weather Rev.* **2006**, *134* (9), 2466–2489.
- Watson, J. G.; Cooper, J. A.; Huntzicker, J. J. The effective variance weighting for least-squares calculations applied to the mass balance receptor model. *Atmos. Environ.* **1984**, *18* (7), 1347–1355.
- Coulter, C. T. *EPA-CMB8.2 Users Manual*; EPA-452/R-04-011; U.S. Environmental Protection Agency: Research Triangle Park, NC, 2004.
- Zheng, M.; Cass, G. R.; Ke, L.; Wang, F.; Schauer, J. J.; Edgerton, E. S.; Russell, A. G. Source apportionment of daily fine particulate matter at Jefferson street, Atlanta, GA, during summer and winter. *J. Air Waste Manage. Assoc.* **2007**, *57* (2), 228–242.
- Zheng, M.; Cass, G. R.; Schauer, J. J.; Edgerton, E. S. Source apportionment of PM<sub>2.5</sub> in the southeastern United States using solvent-extractable organic compounds as tracers. *Environ. Sci. Technol.* **2002**, *36* (11), 2361–2371.
- Zheng, M.; Salmon, L. G.; Schauer, J. J.; Zeng, L. M.; Kiang, C. S.; Zhang, Y. H.; Cass, G. R. Seasonal trends in PM<sub>2.5</sub> source contributions in Beijing, China. *Atmos. Environ.* **2005**, *39* (22), 3967–3976.
- Marmur, A.; Mulholland, J. A.; Russell, A. G. Optimized variable source-profile approach for source apportionment. *Atmos. Environ.* **2007**, *41* (3), 493–505.
- Zielinska, B.; McDonald, J. D.; Hayes, T.; Chow, J. C.; Fujita, E. M.; Watson, J. G. *Northern Front Range Air Quality Study Final Report*; Report to Colorado State University by Desert Research Institute: Reno, NV, 1998.
- Chow, J. C.; Watson, J. G.; Kuhns, H.; Etyemezian, V.; Lowenthal, D. H.; Crow, D.; Kohl, S. D.; Engelbrecht, J. P.; Green, M. C. Source profiles for industrial, mobile, and area sources in the Big Bend Regional Aerosol Visibility and Observational study. *Chemosphere* **2004**, *54* (2), 185–208.
- Cooper, J. A. *Determination of Source Contributions to Fine and Coarse Suspended Particulate Levels in Petersville, AL*; Report to Tennessee Valley Authority by NEA Inc., 1981.



- (26) Paatero, P.; Tapper, U. Positive matrix factorization - a non-negative factor model with optimal utilization of error-estimates of data values. *Environmetrics* **1994**, *5* (2), 111–126.
- (27) Lewis, C. W.; Norris, G. A.; Conner, T. L.; Henry, R. C. Source apportionment of Phoenix PM<sub>2.5</sub> aerosol with the UNMIX receptor model. *J. Air Waste Manage. Assoc.* **2003**, *53* (3), 325–338.
- (28) Norris, G. A.; Vedantham, R.; Wade, K.; Brown, S.; Prouty, J.; Foley, C. *EPA Positive Matrix Factorization (PMF) 3.0 Fundamentals & User Guide*; EPA 600/R-08/1108; U.S. Environmental Protection Agency: Washington, DC, 2008.
- (29) Dennis, R.; Roselle, S.; Gilliam, R.; Arnold, J., High time resolved comparisons for in-depth probing of CMAQ fine particle and gas predictions. In *2004 Models 3 Workshop, Research Triangle Park, NC*, 2004.
- (30) Byun, D. W.; Ching, J. K. S., Eds. *Science Algorithms of the EPA Models-3 Community Multiscale Air Quality (CMAQ) Modeling System*; EPA/600/R-99/030; U.S. Environmental Protection Agency: Research Triangle Park, NC, 1999.
- (31) Houyoux, M. R.; Vukovich, J. M. Updates to the Sparse Matrix Operator Kernel Emissions (SMOKE) Modeling System and Integration with Models-3. The Emission Inventory: Regional Strategies for the Future, Raleigh, NC, October 26–28, 1999; Raleigh, NC, 1999.
- (32) Seaman, N. L. Meteorological modeling for air-quality assessments. *Atmos. Environ.* **2000**, *34* (12–14), 2231–2259.
- (33) Lee, S.; Liu, W.; Wang, Y. H.; Russell, A. G.; Edgerton, E. S. Source apportionment of PM<sub>2.5</sub>: Comparing PMF and CMB results for four ambient monitoring sites in the southeastern United States. *Atmos. Environ.* **2008**, *42* (18), 4126–4137.
- (34) Hansen, D. A.; Edgerton, E. S.; Hartsell, B. E.; Jansen, J. J.; Kandasamy, N.; Hidy, G. M.; Blanchard, C. L. The southeastern aerosol research and characterization study: Part 1-overview. *J. Air Waste Manage. Assoc.* **2003**, *53* (12), 1460–1471.
- (35) Mammur, A.; Unal, A.; Mulholland, J. A.; Russell, A. G. Optimization-based source apportionment of PM<sub>2.5</sub> incorporating gas-to-particle ratios. *Environ. Sci. Technol.* **2005**, *39* (9), 3245–3254.
- (36) Cadle, S. H.; Mulawa, P. A.; Hunsanger, E. C.; Nelson, K.; Ragazzi, R. A.; Barrett, R.; Gallagher, G. L.; Lawson, D. R.; Knapp, K. T.; Snow, R. Composition of light-duty motor vehicle exhaust particulate matter in the Denver, Colorado area. *Environ. Sci. Technol.* **1999**, *33* (14), 2328–2339.
- (37) Hodan, W. *Recommendations for the Update and Improvement of Existing PM<sub>2.5</sub> Split Factors*; Report to EPA by Pacific Environmental Services, Inc.: Research Triangle Park, NC, September 29, 2003.

ES9004703