

**SUPPLEMENTARY MATERIAL:  
EFFICIENT AND EFFECTIVE CALIBRATION OF NUMERICAL MODEL  
OUTPUTS USING HIERARCHICAL DYNAMIC MODELS**

BY YEWEN CHEN <sup>†1</sup>, XIAOHUI CHANG <sup>†2</sup>, BOHAI ZHANG<sup>3</sup> AND HUI HUANG<sup>\*4</sup>

<sup>1</sup>*College of Public Health, University of Georgia, [Yewen.Chen@uga.edu](mailto:Yewen.Chen@uga.edu)*

<sup>2</sup>*College of Business, Oregon State University, [xiaohui.chang@oregonstate.edu](mailto:xiaohui.chang@oregonstate.edu)*

<sup>3</sup>*Guangdong Provincial Key Laboratory of Interdisciplinary Research and Application for Data Science, BNU-HKBU United International College, [bohaizhang@uic.edu.cn](mailto:bohaizhang@uic.edu.cn)*

<sup>4</sup>*Center for Applied Statistics and School of Statistics, Renmin University of China, [huangh89@mail.sysu.edu.cn](mailto:huangh89@mail.sysu.edu.cn)*

In this document, we provide a comprehensive overview of the proposed HDCM, including the used datasets, the competitive models, the selection of the tuning parameters, the simulation study for the proposed VB-spEnKS, the sensitivity of the proposed HDCM with respect to the selection of the tuning parameters, the cross-validation for China's Beijing-Tianjin-Hebei (BTH) data in the summer of 2015, the Diebold-Mariano test, and the Laplace approximation, etc.

**S1. Datasets from multiple sources.** The effectiveness and efficiency of the proposed HDCM are demonstrated using two datasets.

(1) Moderately large datasets. These consist of PM<sub>2.5</sub> concentrations obtained from monitoring stations and the outputs of the Community Multiscale Air Quality (CMAQ) system for China's Beijing-Tianjin-Hebei (BTH) region. The dataset sizes are as follows:  $68 \times 92 = 6,256$  spatio-temporal observations for monitoring stations and  $5,587 \times 92 = 514,004$  raw spatio-temporal outputs for CMAQ, where the number of monitoring stations is 68, the number of CMAQ grid cells is 5,587, and the time length is 92 days.

(2) Large datasets. These include the reanalysis of PM<sub>2.5</sub> outputs of the Nested Air Quality Prediction Modeling System (NAQPMS) and the raw PM<sub>2.5</sub> outputs of the CMAQ system. The dataset sizes are  $6,382 \times 30 = 191,460$  spatio-temporal gridded outputs for NAQPMS and  $16,093 \times 30 = 482,790$  raw gridded outputs for CMAQ, respectively.

**S1.1. Data feature in the BTH region.** As seen in Table S1, the variation of PM<sub>2.5</sub> concentration in most cities of the BTH region is significantly lower in summer than in winter. In general, the calibration for the winter outputs of the numerical model is more challenging than that for the summer outputs.

---

<sup>†</sup>Equal contribution.

\*Corresponding author. Hui Huang, [huangh89@mail.sysu.edu.cn](mailto:huangh89@mail.sysu.edu.cn).

TABLE S1

*Mean, standard deviations (SD), and coefficient of variation (CV) for PM<sub>2.5</sub> concentration predictions ( $\mu\text{g}/\text{m}^3$ ) calculated for 13 cities in the BTH region in both seasons of 2015.*

City	Summer of 2015			Winter of 2015		
	Mean	SD	CV	Mean	SD	CV
Baoding	65.63	32.96	0.50	159.18	111.94	0.70
Beijing	56.18	37.87	0.67	112.80	103.69	0.92
Cangzhou	54.54	25.35	0.46	104.90	78.85	0.75
Chengde	33.49	25.81	0.77	56.38	37.22	0.66
Handan	73.90	28.73	0.39	125.25	90.66	0.72
Hengshui	74.30	28.41	0.38	151.35	113.24	0.75
Langfang	54.74	29.36	0.54	129.38	104.31	0.81
Qinhuangdao	36.36	26.58	0.73	50.81	42.70	0.84
Shijiazhuang	65.81	34.38	0.52	141.47	106.42	0.75
Tangshan	63.17	31.88	0.50	106.37	77.33	0.73
Tianjin	50.18	25.25	0.50	100.05	81.39	0.81
Xingtai	78.08	34.23	0.44	148.89	121.04	0.81
Zhangjiakou	32.45	21.89	0.67	34.40	25.16	0.73
Average	56.83	29.44	0.55	109.33	84.15	0.77

**S1.1.1. Imputing missing data.** In Section 2 of the manuscript, missing data are imputed. Because the missing rate is only 0.28%, this imputation will not affect the data pattern. Of course, some extreme values may be filtered out when we summarize hourly observations to the daily data. However, some of the sites were missing more than 24 hours of hourly observations on some days, and this counts for about 0.28% of our data. To reduce the impact of these missing data on the predictions, we summarize hourly data to daily data and compare the prediction performance of different methods. It is worth noting that the proposed approach is able to directly handle datasets with some hourly data that are completely missing, as it allows imputing the missing data within VB-spEnKS procedure.

Although the proposed approach can impute missing data through VB-spEnKS procedure, missing data is imputed at this stage rather than at the fitting HPCM stage to ensure a fair comparison between different methods by using the same imputed data as the input for all models. In the model comparison, several non-Bayesian methods have difficulties in automatically imputing missing data within their implementation procedures.

**S1.1.2. Data transformation.** In this section, we explore data distributions using histograms (Berrocal, Gelfand and Holland, 2010) and quantile-quantile (Q-Q) boxplots (Rodu and Kafadar, 2022). As seen in Figure S1 (a) and Figure S1 (d), the observed PM<sub>2.5</sub> concentrations show clear right skewness. To stabilize the variance, we carry out two transformations for the observations, namely the square root transformation and logarithmic transformation (Sahu, Gelfand and Holland, 2006). Histograms presented in Figures S1 and Q-Q plots presented in Figures S2 show that the square root transformed data are closer to the normal distribution compared to the log-transformed data. Therefore, in this work, we use the square-root transformation.

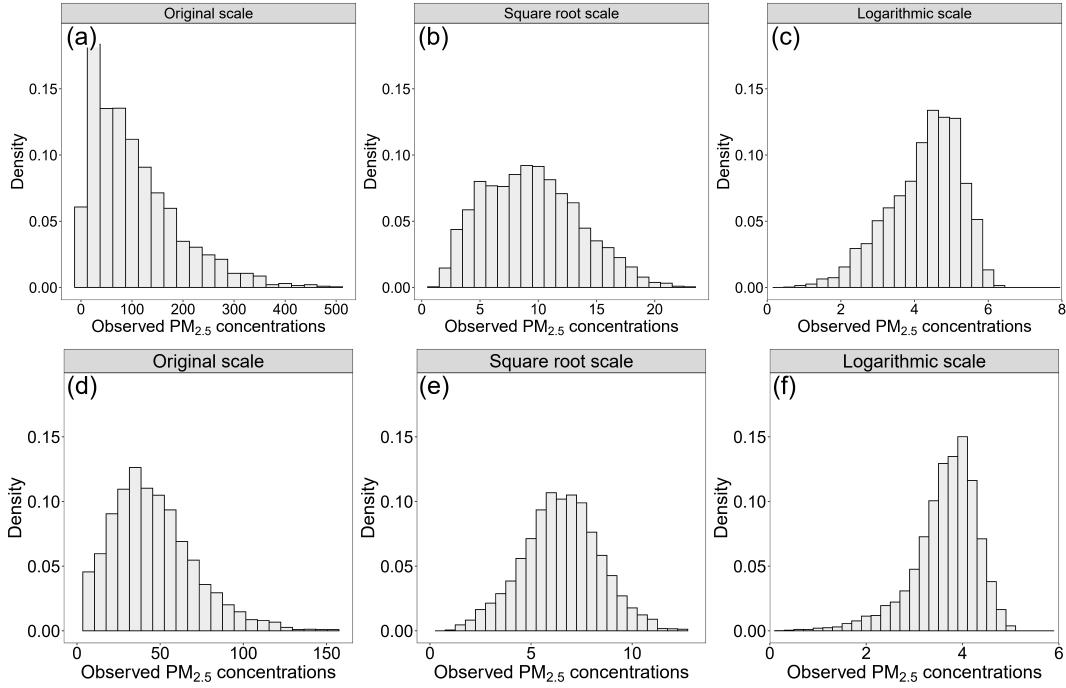


FIG S1. Histograms of observed PM<sub>2.5</sub> concentrations at three scales, including the original, square root, and logarithmic scales. (a)-(c) In the winter of 2015. (d)-(f) In the summer of 2015.

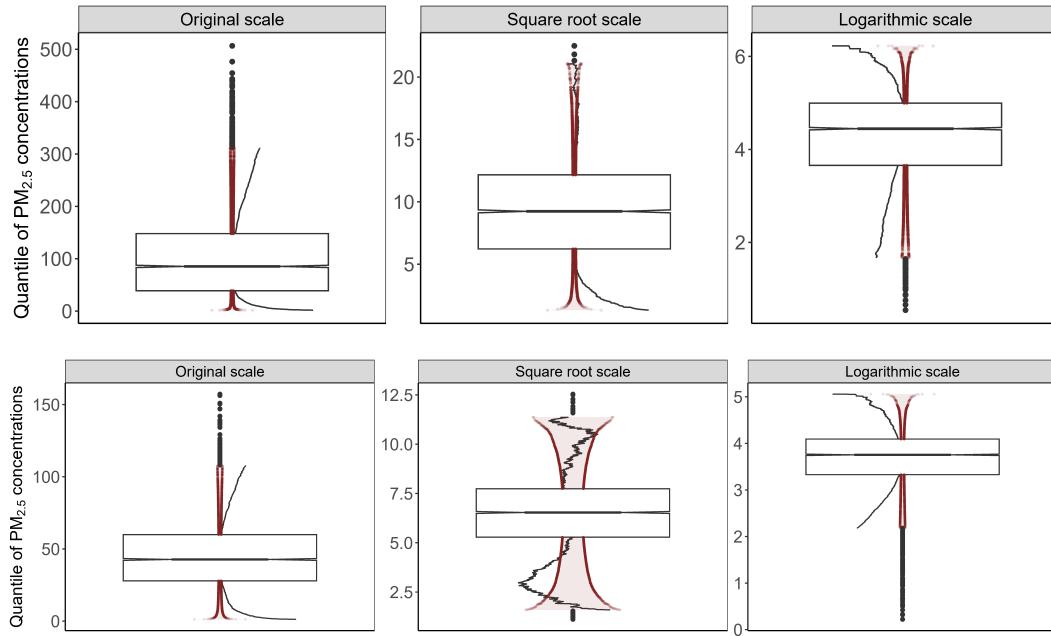


FIG S2. Q-Q boxplots of observed PM<sub>2.5</sub> concentrations at three scales, including the original, square root, and logarithmic scales. (a)-(c) In the winter of 2015. (d)-(f) In the summer of 2015.

**S1.2. Large datasets.** Figure S3 displays spatial distributions of the grid cells under different gridding systems, namely CMAQ and NAQPMS. Both datasets cover an area that is much larger than the BTH region.

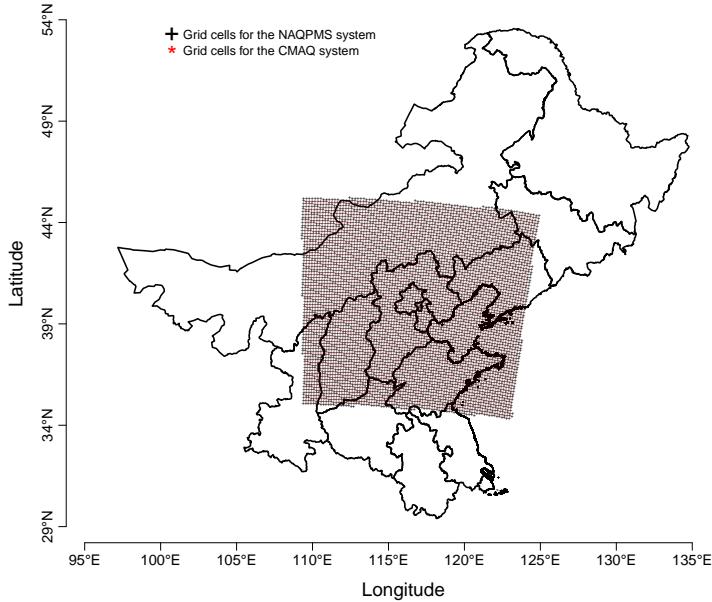


FIG S3. Maps of grid cells of CMAQ and NAQPMS. The symbols “+” represent the centroids of 16,093 9km CMAQ grids. The symbols “\*” denote the centroids of 6,382 15km NAQPMS grids.

As seen in Section 6 of the manuscript, both the proposed HDCM and STAR require triangulated meshes, and an equal area partition method is applied for HDCM. As an illustration, we present a triangulated mesh with triangle vertices of  $m = 10,103$  and subregions of  $R = 9$  in Figure S4.

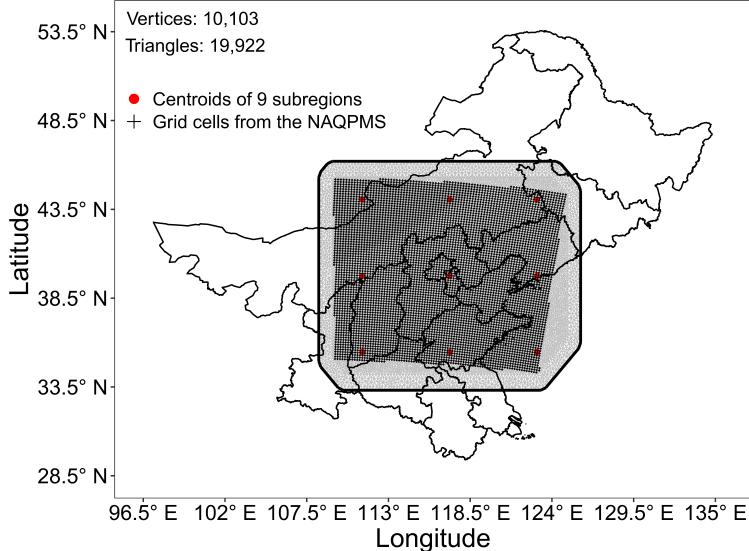


FIG S4. Triangulated mesh for the reanalysis dataset with 9 subregions, whose centroids are marked with red dots.

**S2. Competing models.** In this section, we describe how several competing models are implemented.

**S2.1. Spatially-varying coefficients (SVC) model.** The following version of SVC is used for the model comparison:

$$(S1) \quad y_t(s) = \mu_t(s) + w_t(s) + \varepsilon_t(s).$$

where the zero-mean spatial process  $w_t(s)$  and the zero-mean white noise  $\varepsilon_t(s)$  are mutually and serially independent.

This SVC model differs from the proposed HDCM model in the mean term  $\mu_t(s)$  and the random effects term  $w_t(s)$ . For  $\mu_t(s)$ , we employ the weighted CMAQ PM<sub>2.5</sub> output as a single covariate to improve the performance of SVC model, i.e.,  $\mu_t(s) = \beta_0 + \beta_1 x_{t,1}(C_s)$  where  $x_{t,1}(C_s)$  represents the square root of the weighted PM<sub>2.5</sub> output of CMAQ. This is similar to Berrocal et al. (2020) where SVC uses CMAQ output as a single covariate. Using SVC with the single covariate for the BTH data in the winter of 2015 improves the performance by 11.80% in average RMSE and 11.34% in average CRPS based on the leave-one-city-out cross-validation (LOCOCV) of Section 4 of the manuscript.

In addition, at each time  $t$ , SVC (S1) models  $w_t(s)$  as spatial stationary Gaussian random processes with mean  $\mathbf{0}$  and an exponential covariance function,  $\text{Cov}(w_t(s_i), w_t(s_j)) = \tau_w^2 \exp(-\|s_i - s_j\|/\phi_w)$ . The unknown parameters including  $\tau_w^2$ ,  $\phi_w$ , and others appearing in (S1) can be estimated using the Markov chain Monte Carlo (MCMC) algorithm. We use the `spBayes` package (<https://cran.r-project.org/web/packages/spBayes/>) to fit SVC model.

**S2.2. Universal kriging (UK).** UK follows closely the model specification of SVC, where the prediction performance of UK is also improved by using the square root of the weighted PM<sub>2.5</sub> output of CMAQ as the single covariable. The difference is that parameters

other than the covariate coefficient vector  $\beta = (\beta_0, \beta_1)$  in the term  $\mu_t(s) = \beta_0 + \beta_1 x_{t,1}(C_s)$  are first estimated using the restricted maximum likelihood method (details can be referred to the `geoR` package documentation <https://cran.r-project.org/web/packages/geoR/>). We then estimate  $\beta$  and implement UK interpolation using the `gstat` package (<https://cran.r-project.org/web/packages/gstat/>).

**S2.3. Spatiotemporally-varying coefficient (STVC) model.** Following Section 3.2 of Berrocal, Gelfand and Holland (2010), the spatial SVC (S1) can be extended to a spatio-temporal downscaler model by allowing spatiotemporally varying coefficients (STVC), where both the intercept and the slope of CMAQ variable vary with respect to time. STVC is given by

$$(S2) \quad \begin{cases} y_t(\mathbf{s}) = \beta_0 + \beta_{0,t} + (\beta_1 + \beta_{1,t})x_{t,1}(C_s) + \sum_{k=2}^4 \beta_k x_{t,k}(C_s) + w_t(\mathbf{s}) + \varepsilon_t(\mathbf{s}) \\ \beta_{l,t} = \rho_l \beta_{l,t-1} + \eta_{l,t}, \text{ for } l = 0, 1, \end{cases}$$

where  $x_{t,1}(C_s), \dots, x_{t,4}(C_s)$  represent the square root of the weighted average CMAQ PM<sub>2.5</sub> output, the weighted average surface temperature, the weighted average surface pressure, the weighted average northern and eastern cumulative wind powers at location  $\mathbf{s}$  and time point  $t$ , respectively. The random effects  $\beta_{l,t}$  and the spatial process  $w_t(\mathbf{s})$  are assumed to be mutually and serially independent, the same for  $\beta_{l,t}$  and  $\varepsilon_t(\mathbf{s})$ , for  $w_t(\mathbf{s})$  and  $\varepsilon_t(\mathbf{s})$ , for  $\beta_{l,t-1}$  and  $\eta_{l,t}$ , and for  $\eta_{l,t}$  and  $\varepsilon_t(\mathbf{s})$ .

The definition of  $w_t(\mathbf{s})$  is the same to SVC (S1), and  $\boldsymbol{\eta}_{l,t} = (\eta_{l,1}, \dots, \eta_{l,N_t})^T \sim \mathcal{N}(0, \delta_l^2 \mathbf{I}_{N_t})$ . Especially, the performance of (S2) for the BTH data can be improved by assuming  $\beta_{1,t} = 0$  based on LOCOCV, and thus the setting of  $\beta_{1,t} = 0$  is used in STVC. STVC is estimated based on MCMC algorithms which are available through the R package `sptTDyn` (Bakar, Kokic and Jin, 2016).

**S2.4. First-order spatiotemporal autoregression (STAR) model.** We follow Section 7.2 of Blangiardo and Cameletti (2015) and use STAR model of the form:

$$(S3) \quad \begin{cases} y(\mathbf{s}) = \beta_0 + \sum_{k=1}^4 \beta_k x_{t,k}(C_s) + w_t(\mathbf{s}) + \varepsilon_t(\mathbf{s}) \\ w_t(\mathbf{s}) = \rho w_{t-1}(\mathbf{s}) + \eta_t(\mathbf{s}), \end{cases}$$

where the zero-mean spatio-temporal process  $w_t(\mathbf{s})$  and the zero-mean white noise  $\varepsilon_t(\mathbf{s})$  are mutually and serially independent, the same for  $w_{t-1}(\mathbf{s})$  and  $\eta_t(\mathbf{s})$ , and for  $\eta_t(\mathbf{s})$  and  $\varepsilon_t(\mathbf{s})$ .

The initial Gaussian process  $w_0(\mathbf{s}) = (w_t(s_1), \dots, w_t(s_n))^T \sim \mathcal{N}(\boldsymbol{\mu}_0, \boldsymbol{\Sigma}_0)$ , and the Gaussian random process  $\boldsymbol{\eta}_t = (\eta_t(s_1), \dots, \eta_t(s_n))^T \sim \mathcal{N}(\mathbf{0}, \boldsymbol{\Sigma}_\eta)$ . In this work, STAR is implemented through the R package `INLA` (Lindgren and Rue, 2015), where a Gaussian process is linked to a Gaussian Markov random field (GMRF, Rue and Held (2005)) on a triangulated mesh via a stochastic partial differential equation (Lindgren, Rue and Lindström, 2011), and the parameters are then estimated by the integrated nested Laplace approximation (INLA, Rue, Martino and Chopin (2009)). This a powerful Bayesian approach for handling large spatio-temporal datasets (Rue et al., 2017).

**S2.5. Random forest (RF).** For the random forest model, we estimate an unknown function  $f(\cdot)$  by minimizing the out-of-bag error:

$$\sum_{t=1}^T \sum_{s=1}^n \text{Loss}\left\{y_t(\mathbf{s}), f(x_{t,1}(C_s), \dots, x_{t,4}(C_s))\right\}.$$

The random forest recursively partitions the covariate space to build the trees and uses the weighted average of the predictions across all trees as the final prediction. In this work, we use 500 trees to build the predictive model that is implemented through the `rfinterval` package (<https://github.com/haozhestat/rfinterval>).

### S3. HDCM.

**S3.1. Mapping matrix for HDCM and precision matrices of IDE models.** In [Lindgren, Rue and Lindström \(2011\)](#), the construction of the matrix  $\mathbf{G}$  depends on a piecewise linear basis  $\{\psi_k(s)\}_{k=1}^m$  on the mesh, and  $G_{k,k'} = \langle \nabla \psi_k \nabla \psi_{k'} \rangle$ , where  $G_{k,k'}$  is the  $(k, k')$ th element of  $\mathbf{G}$ ,  $\langle \cdot \rangle$  is defined in Equation (7) of [Lindgren, Rue and Lindström \(2011\)](#), and the vertex index  $k, k' = 1, \dots, m$ . In our work,  $h(i, k)$  plays an equivalent role to  $\psi_k(s_i)$  in their method, where  $h(i, k)$  is the  $(i, k)$ th element of  $\mathbf{H}$  and is associated with the location  $s_i$  and the vertex  $B_k$ . We did not construct our matrix  $\mathbf{G}$  as a function of the triangle geometry using the definition of  $G_{k,k'}$  described earlier.

This is because the piecewise linear assumption of the bases is not suitable for our data based on additional numerical experiments, which may be related to the distribution of data. Especially, because the BTH data are sparsely observed in space and clustered in urban areas, it would be difficult to identify  $v_t$  located in rural areas with no stations using piecewise linear bases, and the cross-validation also shows that the piece-wise linear basis for the BTH data does not work well.

Moreover, we require the matrix  $\mathbf{G}$  to be sparse. Under the setting of [Lindgren, Rue and Lindström \(2011\)](#) (i.e.,  $G_{k,k'} = \langle \nabla \psi_k \nabla \psi_{k'} \rangle$ ), it is usually difficult to ensure the sparsity of  $\mathbf{G}$  when the piecewise linear basis  $\psi$  is replaced by other functions such as the Wendland function used in this work. Therefore, in this work, we employ a Laplacian matrix to define  $\mathbf{G}$  and specify the corresponding precision matrix as a function of  $\mathbf{G}$  as seen in Section 3.2.2 of the manuscript; a similar definition for the precision matrix of Gaussian random field is also adopted in [Bolin, Wallin and Lindgren \(2019\)](#).

**S3.2. Prior distributions for HDCM.** The prior distributions used in this work are provided in Table S2.

TABLE S2  
*Prior distributions of parameters in  $\Theta$ .*

Parameter	Prior distribution
$\beta$	$\mathcal{N}(\mathbf{0}, 10^5 \mathbf{I}_5)$
$\sigma^2$	IG(2, 1)
$\tau^{2,(r)}$	Gamma(2, 1)
$\tau_0^{2,(r)}$	Gamma(2, 1)
$\zeta^{2,(r)}$	Uniform(0, 400)
$\zeta_0^{2,(r)}$	Uniform(0, 400)
$\theta_1^{(r)}$	$\mathcal{N}(0.0001, 10^5)$
$\theta_2^{(r)}$	Uniform(0.001, 1)

For all  $r = 1, \dots, R$ .

**S4. Selection of tuning parameters.** Our approach needs to select four tuning parameters: (1) the number of ensembles  $N_e$ , (2)  $c_h \in (0, 1]$  related to the mapping matrix  $\mathbf{H}$ , (3)  $c_s \in (0, 1]$  for the covariance tapering in space, and (4)  $c_t \in \{1, 2, \dots, N_t\}$  for the covariance tapering in time. In this work, we adopt the commonly used ensemble size in the literature ([Mitchell, Houtekamer and Pellerin, 2002](#); [Houtekamer and Zhang, 2016](#)) and let  $N_e = 100$ . The other three parameters are obtained through LOOCV of Section 4 of the manuscript.

**S4.1. Cross-validation for tuning parameters.** In LOCOCV, considering  $J$  combinations for given  $c_h$ ,  $c_s$ , and  $c_t$ , an average RMSE across all 13 cities can be calculated with the predictions for each of the combinations, denoted by  $\text{avgRMSE}_j$ , where  $j = 1, 2, \dots, J_0$ . Assuming  $J_0 = \min_j \{\text{avgRMSE}_j\}$ , then the  $J_0$ th combination of tuning parameters is specified as the choice of these three parameters.

As seen in [Kirchgessner, Nerger and Bunse-Gerstner \(2014\)](#), the tuning parameters usually depend on the distributions of observations, and the parameters themselves could also be interrelated.

In the setting of HDCM,  $c_h$  and  $c_t$ , in general, depend on the spatial range and temporal range, respectively, which can be utilized to facilitate the grid search in the cross-validation. The following strategies are adopted in LOCOCV:

- a. For  $c_h \in (0, 1]$ , the length of this interval can be further reduced to improve the grid search in cross-validation by covariograms (or variograms). Based on covariograms, for a given  $c_t$ , the interval  $(0, 1]$  can be reduced to  $[c_h^{\min}, c_h^{\max}]$  by setting  $c_h^{\min} d_{\max}^H = \phi_s^{(1)}$  and  $c_h^{\max} d_{\max}^H = \phi_s^{(2)}$ , where  $d_{\max}^H$  is defined in Section 3.2.1 of the manuscript,  $\phi_s^{(1)}$  and  $\phi_s^{(2)}$  correspond to the distances where the spatial correlation coefficients are equal to 0.9 and 0.05, respectively. Because  $0 < \phi_s^{(1)} \leq d_{\max}^H$  and  $0 < \phi_s^{(2)} \leq d_{\max}^H$ , the length of  $[c_h^{\min}, c_h^{\max}] = [\phi_s^{(1)}/d_{\max}^H, \phi_s^{(2)}/d_{\max}^H]$  is smaller than that of  $(0, 1]$ . In practice, although these two values of  $\phi_s^{(1)}$  and  $\phi_s^{(2)}$  are usually not directly obtained from covariograms, we can fit a parametric spatio-temporal covariance model to an empirical covariogram such as Figure 2 of the manuscript and then estimate them using the fitted parametric covariance model. In the data analysis, we used  $c_h \in [0.1, 0.31]$  for the BTH data and  $c_h \in (0.03, 0.17]$  for the reanalysis data.
- b. For  $c_t$ , the empirical covariograms can be also used to narrow down the grid search in LOCOCV. For example, Figure 2 shows the temporal range is likely to be less than 5 days, and we can select a  $c_t$  value from  $\{1, 2, \dots, 5\}$  other than  $\{1, 2, \dots, N_t\}$ , where  $N_t$  is the length of time points.
- c. For the covariance tapering parameter in space  $c_s$ , [Kirchgessner, Nerger and Bunse-Gerstner \(2014\)](#) showed  $c_s = 8\sqrt{\frac{N_e}{40}}dx$  in dense observations from regular grids, where  $dx$  represents the grid spacing. Although the formula does not provide an optimal value of  $c_s$  for the triangle mesh in the proposed HDCM, we can use this definition to construct an interval such that  $c_s \in [8\sqrt{\frac{N_e}{40}}dx_{(1)}/d_{\max}^B, 8\sqrt{\frac{N_e}{40}}dx_{(2)}/d_{\max}^B]$ , where  $d_{\max}^B = \max_{l,l'} \{\|B_l^{(r)} - B_{l'}^{(r)}\|\}$ . The length of this interval is smaller than that of  $(0, 1]$ . Based on numerical experiments, the predictive performance of HDCM is rather competitive when  $dx_{(1)}$  is set as 2.5% quantile of  $\{\|B_l^{(r)} - B_{l'}^{(r)}\|\}_{r,r'=1}^m$ , and  $dx_{(1)}$  as 95% quantile of  $\{\|B_l^{(r)} - B_{l'}^{(r)}\|\}_{l,l'=1}^m$ .

**S4.2. Sensitivity analysis with respect to tuning parameters.** In this section, we investigate the proposed HDCM with respect to the sensitivity of the selection of several tuning parameters. The sensitivity analysis is carried out based on the cross-validation procedure described above (i.e., LOCOCV). Because the sensitivity analysis results for all 13 cities are similar, here we choose Zhangjiakou as an illustration. Specifically, we consider all monitoring data from Zhangjiakou in the winter as the test set and all data from the remaining twelve cities in the winter as the training set. We evaluate the predictive performance of HDCM in terms of RMSE with respect to various tuning parameters.

Although the tuning parameter  $c_h$  can be selected from the interval  $(0, 1]$ , we found that the performance of model is better when  $c_h \in (0, 0.5]$  than when  $c_h \in (0.5, 1]$ . Therefore, we consider  $c_h = \{0.1, 0.15, 0.2, 0.25, 0.3, 0.35, 0.4, 0.45, 0.5\}$  in the sensitivity analysis. We also set the two covariance tapering parameters to  $c_s = \{0.1, 0.2, 0.3, 0.4, 0.5, 0.6, 0.7, 0.8, 0.9, 1\}$  and  $c_t = \{1, 2, 3, 4, 5\}$ . The size of ensemble members  $N_e = \{30, 50, 70, 100\}$ .

RMSEs for various combinations of the tuning parameters for the training data and test data are presented in Figure S5 and Figure S6, respectively. Below are the results:

- (a) For  $c_h$ , RMSE increases as  $c_h$  gets closer to 0.5 while the predictions are robust when  $c_h \in (0.1, 0.3]$ .
- (b) For  $c_s$ , when  $c_h < 0.3$ , the predictions with  $c_s < 0.5$  in the test dataset are more accurate than those with  $c_s > 0.5$  while the predictions in the training dataset are robust for all  $c_s$ .
- (c) For  $c_t$ , the predictions in the training dataset are usually robust in almost all cases while the prediction performance in the test dataset is better for a smaller  $c_t$ . This suggests  $c_t = 1$  is an optimal choice.
- (d) For  $N_e$ , the predictions in almost all cases tend to be more accurate for a larger  $N_e$ .

In addition, the computation becomes more expensive when  $c_s$  is closer to 1 and/or  $c_t$  is closer to 5. More specifically, the computation time when  $c_s = 1$  and  $c_t = 5$  is nearly 3 times of that when  $c_s = 0.1$  and  $c_t = 1$ , while these two parameters are increased by a factor of 10 and 5, respectively; see Figure S7. In summary, under the current setting, the computational cost of the proposed approach is almost linear with respect to these two parameters.

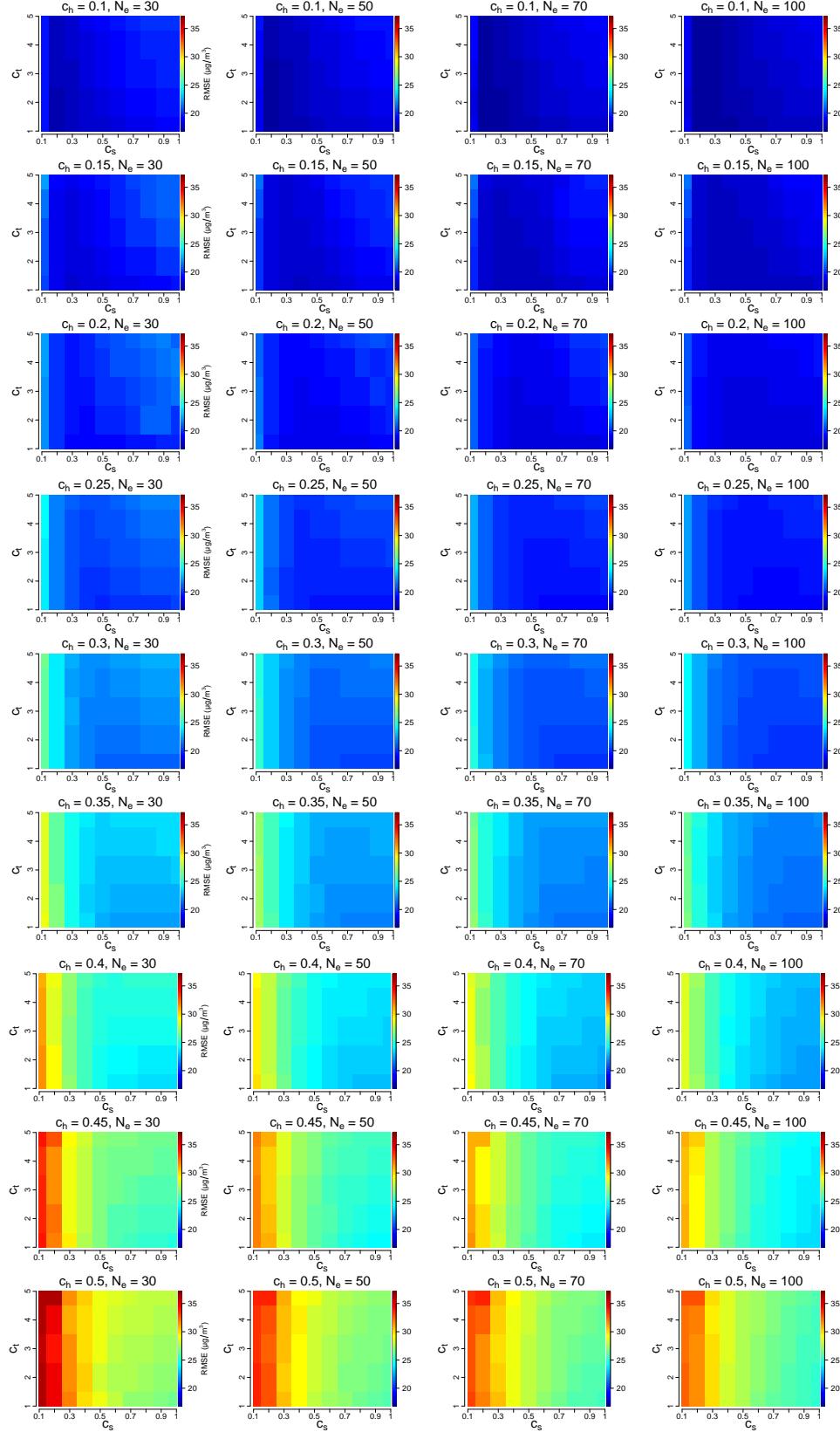


FIG S5. Root mean squared error (RMSE) of the training dataset for various tuning parameters.

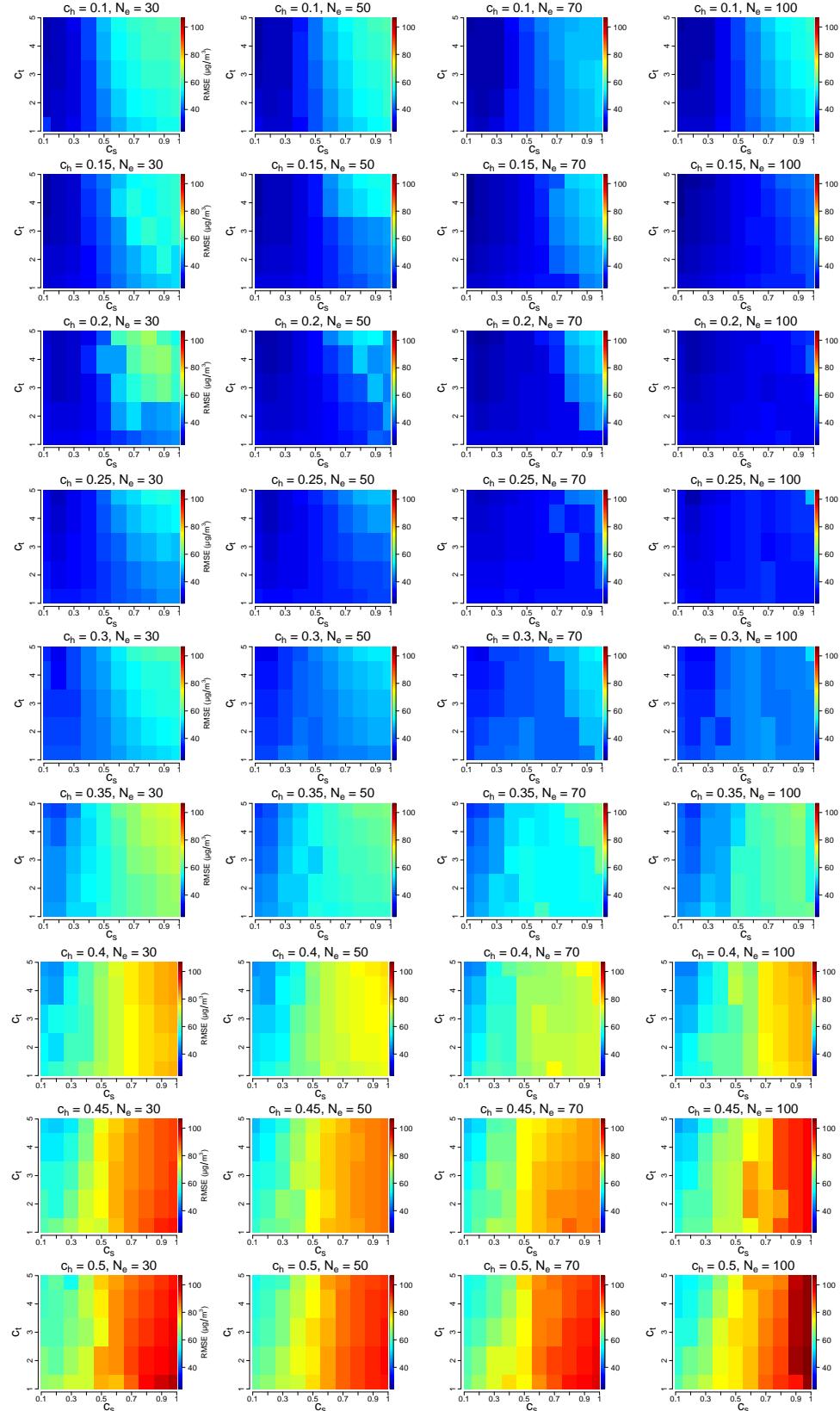


FIG S6. Root mean squared error (RMSE) of the test dataset for various tuning parameters.

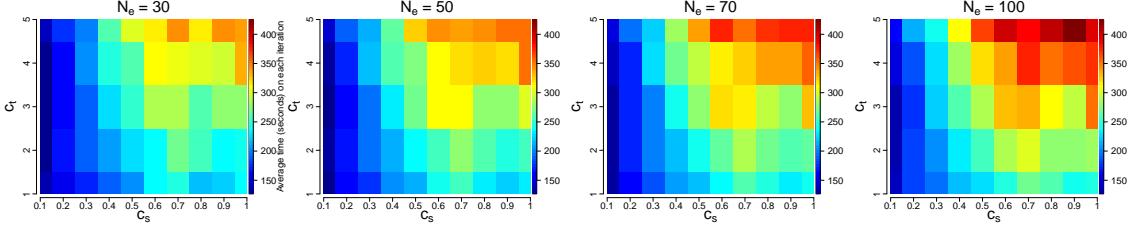


FIG S7. Average running time (seconds) for each iteration.

**S5. Comparing VB-spEnKS with MCMC-KS.** VB and spEnKS (VB-spEnKS) is the approximation of MCMC and KS (MCMC-KS). We first apply both VB-spEnKS and MCMC-KS to HDCM on the simulated data, and then compare their performances.

**S5.1. Simulation design.** We set the spatial domain  $\mathcal{D}'$  to be the BTH region with 2,499 equally spaced grids, which are the regular CMAQ grids as shown in Figure 1 (b) of the manuscript and denoted as  $\{C_{i''}\}_{i''=1}^{2,499}$ . The time domain  $\mathcal{T}'$  is set to be  $[0, 1]$  with 20 equally spaced time points. For any  $s \in \mathcal{D}'$  at time point  $t \in \mathcal{T}'$ , we use the following model to generate data,

$$(S4) \quad y_t(s) = \beta_0 + \beta_1 x_{1,t}(s) + \beta_2 x_{2,t}(s) + w_t(s) + \varepsilon_t(s),$$

where the covariate  $x_{1,t}(s)$  is generated from the standard normal distribution,  $x_{2,t}(s)$  is from a spatial Gaussian random field with mean  $\mathbf{0}$  and an isotropic exponential covariance kernel, i.e., for location  $s$  and  $s'$ ,  $\text{Cov}(x_{2,t}(s), x_{2,t}(s')) = \exp(-10\|s - s'\|/d_{\max}^s)$  with  $d_{\max}^s = \max_{s,s'} \{\|s - s'\|\}$ , the spatio-temporal random process  $w_t(s)$  is mutually and serially independent of the error term  $\varepsilon_t(s)$ , and the error term  $\varepsilon_t(s) \sim N(0, \sigma^2)$ . The random process  $w_t(s)$  follows a zero-mean stationary space-time Gaussian process with a kernel specified in the Gneiting class (Gneiting, 2002), i.e., for any  $(s, t)$  and  $(s', t')$ , the covariance between  $w_t(s)$  and  $w_{t'}(s')$  is given by

$$(S5) \quad \text{cov}(w_t(s), w_{t'}(s')) = \sigma_w^2 (1 + \Psi_t(\tau))^{-\delta/2} \Psi_s(u/\sqrt{\Psi_t(\tau)}),$$

where  $\tau = |t - t'|$ ,  $\sigma_w > 0$ ,  $\delta \geq 2$ ,  $\Psi_t(\tau) = (|\tau/\phi_t|^a + 1)^{-b/a}$  for  $\phi_t > 0$ ,  $a \in (0, 2]$  and  $b > 0$ ,  $u = \|s - s'\|$ ,  $\Psi_s(u) = 2^{1-\nu} \Gamma^{-1}(\nu) (\sqrt{2\nu} u / \phi_s)^\nu K_\nu(\sqrt{2\nu} u / \phi_s)$  for  $r \geq 0$ ,  $\phi_s > 0$ ,  $\nu > 0$ ,  $\Gamma(\cdot)$  is the Gamma function, and  $K_\nu$  is the modified Bessel function of the second kind.

We set model parameters of (S4) to  $\beta = (15, 1, 1)$ ,  $\sigma_w^2 = 1$ ,  $\delta = 2$ ,  $a = 1$ ,  $b = 1.5$ ,  $\phi_t = 0.2$  for  $\Psi_t$ , and  $\nu = 3$ ,  $\phi_s = 100$  for  $\Psi_s$ .

The simulation consists of the following steps:

- Step 1: Sample  $w_t(s)$  using (S5), where  $s \in \mathcal{S} = \{C_1, \dots, C_{2499}\}$ , and  $t = 0, 1/19, \dots, 1$ .
- Step 2: Generate synthetic data  $y_t(s)$  using (S4).
- Step 3: Sample  $n$  spatial points from  $\mathcal{S}$  and denote as  $\mathcal{S}^{(1)} = \{s'_1, \dots, s'_n\}$ . Let  $\mathcal{S}^{(2)} = \mathcal{S} \setminus \mathcal{S}^{(1)}$ . The training set is from  $\mathcal{S}^{(1)}$  while the test set is from  $\mathcal{S}^{(2)}$ .
- Step 4: Fit HDCM using the training set through VB-spEnKS and MCMC-KS algorithms, and generate predictions of the test data  $y_t(s)$  in  $\mathcal{S}^{(2)}$ , denoted by  $\hat{y}_t(s)$ .
- Step 5: Repeat Steps 2 – 4 to assess the estimates of the parameters (e.g.,  $\beta$ ), the predictions  $\hat{y}_t(s)$ , and the recovered surface for the random process  $w_t(s)$ .

In step 4, we set four tuning parameters to  $c_h = 0.23$ ,  $c_s = 0.3$ ,  $c_t = 1$ , and  $N_e = 100$ . In step 5, the parameter estimates are evaluated using standard deviation (SD) and mean squared error (MSE). Under each setting, we repeat the procedure 40 times and report the results in Table S3.

**S5.2. Simulation results.** Table S3 shows the parameter estimates and prediction performance of VB-spEnKS are very similar to that of MCMC-KS. VB-spEnKS in some cases produces even smaller MSEs than MCMC-KS, indicating that VB-spEnKS can approximate MCMC-KS very well. Moreover, VB-spEnKS has higher computational efficiency than MCMC-KS as it took MCMC-KS more than two weeks to complete the 40 independent simulations but less than 5 hours for the same task on the same workstation. On average, VB-spEnKS converges after 15 iterations. In Table S3, VB-spEnKS results are the summaries after about 15 iterations and MCMC-KS results are from 2,500 samples after a burn-in of 2,500. All these results demonstrate that VB-spEnKS can provide a competitive performance within dozens of iterations.

TABLE S3

The estimates, standard deviations (SD), and mean squared errors (MSEs) of the regression coefficients  $\beta$  and the nugget effect  $\sigma^2$  when the training spatial sample size  $n = 100$ . The results are based on 40 independent simulations and obtained by fitting HDCM model with two algorithms: (1) Markov chain Monte Carlo and Kalman smoother (MCMC-KS), (2) Variational Bayes and spatial-partitioning-based Ensemble Kalman smoother (VB-spEnKS) with a subregion.

$n$	Algorithm	Criteria	$\beta_0 (= 15)$	$\beta_1 (= 1)$	$\beta_2 (= 1)$	$\sigma^2 (= 0.1)$	Prediction
100	MCMC-KS	Estimate	14.8558	0.9994	0.9997	0.1235	-
		SD	0.0946	0.0095	0.0166	0.0055	-
		MSE	0.0297	0.0001	0.0003	0.0006	0.1678
	VB-spEnKS	Estimate	14.9769	0.9997	0.9998	0.1452	-
		SD	0.0726	0.0094	0.0169	0.0061	-
		MSE	0.0059	0.0001	0.0003	0.0021	0.1866

In addition, our proposed model (2) for  $w_t(s)$  in the manuscript successfully captures most of the spatio-temporal patterns in the true spatio-temporal process. As an illustration, Figure S8 demonstrates the excellent performance of HDCM in recovering the true  $w_t(s)$  from a simulated diffusion process for  $n = 100$ .

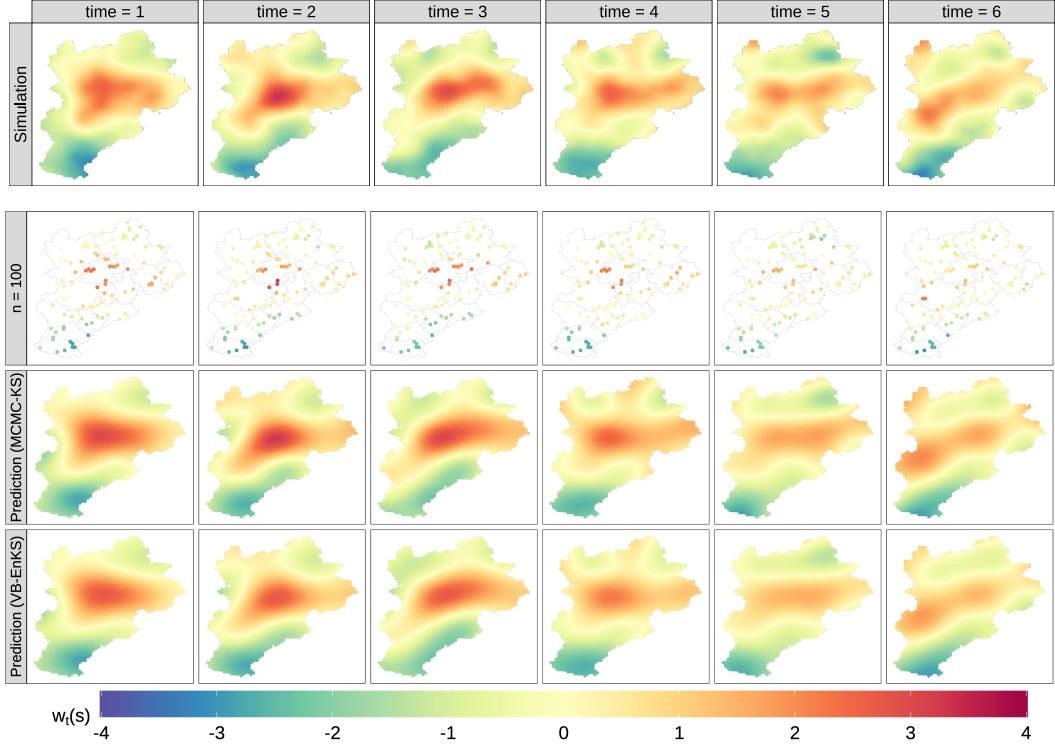


FIG S8. Maps of the random process  $w_t(s)$  in the Gneiting space-time covariance model. From top to bottom: simulated  $w_t(s)$  map, map of 100 sampling locations, predicted  $w_t(s)$  map using MCMC-KS based on the 100 locations, and predicted  $w_t(s)$  map using VB-spEnKS based on the 100 locations. The time horizon is from time = 15 to time = 20.

## S6. Cross validation.

**S6.1. Model comparisons for the summer data.** In this section, we consider the data with the square root scale in the summer of 2015, and the results of LOCOCV are reported in Table S4 and Table S5.

TABLE S4

*Averaged RMSE (Root Mean Squared Error) and CRPS (Continuous Rank Probability Score) for PM<sub>2.5</sub> concentration predictions ( $\mu\text{g}/\text{m}^3$ ) calculated for 13 cities in the BTH region using leave-one-city-out cross-validation from seven models: (1) CMAQ numerical model output, (2) universal kriging (UK), (3) random forest model (RF), (4) spatially-varying coefficient downscaling model (SVC), (5) spatiotemporally-varying coefficient downscaling model (STVC), (6) first-order spatio-temporal autoregression (STAR), and (7) the proposed hierarchical dynamic calibration model (HDCM). The smallest RMSE and CRPS are in bold, and the second smallest ones are underlined. Differences between methods at a 5% significance level are detailed in Table S8. Daily data from June 1, 2015, and August 31, 2015, are considered.*

City	RMSE						
	CMAQ	UK	RF	SVC	STVC	STAR	HDCM
Baoding	36.08	24.92	25.41	<b>22.17</b>	23.96	22.96	<u>22.74</u>
Beijing	36.87	<u>22.00</u>	27.17	<b>21.57</b>	24.19	22.77	22.50
Cangzhou	34.41	20.92	20.32	<u>17.70</u>	18.42	<b>16.63</b>	19.07
Chengde	25.04	18.74	19.53	<b>16.68</b>	18.12	<u>17.59</u>	18.82
Handan	34.33	27.21	23.33	<u>19.12</u>	26.79	20.12	<b>19.10</b>
Hengshui	44.14	30.53	25.37	<u>21.18</u>	24.69	<b>20.47</b>	21.19
Langfang	27.52	15.48	17.56	<b>13.28</b>	16.08	<u>13.61</u>	13.78
Qinhuangdao	24.16	28.11	23.47	24.44	<u>23.20</u>	23.81	<b>20.93</b>
Shijiazhuang	34.23	27.57	28.43	<u>25.90</u>	28.44	<b>25.24</b>	26.23
Tangshan	43.65	36.51	<u>25.39</u>	27.61	25.50	<b>24.42</b>	25.65
Tianjin	24.22	15.04	17.50	<b>14.01</b>	18.68	<u>14.30</u>	14.69
Xingtai	40.02	28.73	26.48	<b>20.77</b>	28.20	21.87	<u>21.09</u>
Zhangjiakou	23.98	23.81	19.35	27.39	<u>17.53</u>	23.57	<b>16.37</b>
Average	32.97	24.58	23.02	20.91	22.60	<u>20.57</u>	<b>20.17</b>

City	CRPS						
	CMAQ	UK	RF	SVC	STVC	STAR	HDCM
Baoding	21.27	13.19	13.43	<b>11.97</b>	12.52	12.06	<u>12.02</u>
Beijing	20.51	<u>11.87</u>	14.48	<b>11.76</b>	12.93	12.24	12.10
Cangzhou	23.13	11.49	11.26	<u>9.78</u>	10.67	<b>9.44</b>	10.44
Chengde	12.56	10.41	10.10	<u>9.51</u>	<b>9.37</b>	9.60	9.61
Handan	20.29	15.37	12.99	<u>11.34</u>	15.01	11.39	<b>11.00</b>
Hengshui	30.35	17.85	13.90	11.80	13.74	<b>11.27</b>	<u>11.79</u>
Langfang	15.94	9.35	10.32	<u>9.29</u>	9.53	<b>8.91</b>	9.61
Qinhuangdao	12.77	15.66	12.82	13.39	<u>12.62</u>	13.02	<b>11.46</b>
Shijiazhuang	19.27	14.94	15.69	14.46	15.50	<b>13.67</b>	<u>14.38</u>
Tangshan	25.75	20.75	<u>14.43</u>	15.61	14.64	<b>13.98</b>	14.72
Tianjin	13.75	8.81	9.99	<b>8.55</b>	10.70	<u>8.61</u>	8.84
Xingtai	23.55	16.09	14.53	<b>11.62</b>	15.56	12.01	<u>11.96</u>
Zhangjiakou	13.32	13.28	10.61	15.19	<u>9.74</u>	13.19	<b>9.06</b>
Average	19.42	13.77	12.66	11.86	12.50	<u>11.49</u>	<b>11.30</b>

TABLE S5

*Averaged MAE (Mean Absolute Error) and FAC2 (Fraction of Predictions within a factor of 2) for PM<sub>2.5</sub> concentration predictions ( $\mu\text{g}/\text{m}^3$ ) calculated for 13 cities in the BTH region using leave-one-city-out cross-validation from eight models: (1) CMAQ numerical model output, (2) universal kriging (UK), (3) random forest model (RF), (4) spatially-varying coefficient downscaling model (SVC), (5) spatiotemporally-varying coefficient downscaling model (STVC), (6) first-order spatio-temporal autoregression (STAR), and (7) the proposed hierarchical dynamic calibration model (HDCM). The smallest MAE and FAC2 are in bold, and the second smallest ones are underlined. Daily data from June 1, 2015, and August 31, 2015, are considered.*

City	MAE						
	CMAQ	UK	RF	SVC	STVC	STAR	HDCM
Baoding	27.76	17.86	18.11	<b>15.77</b>	16.71	<u>16.28</u>	16.37
Beijing	26.37	15.38	19.49	<b>14.84</b>	17.44	<u>16.14</u>	<u>15.19</u>
Cangzhou	28.16	15.59	15.615	<u>13.10</u>	14.26	<b>12.53</b>	14.22
Chengde	15.83	13.71	13.23	<b>11.74</b>	<u>12.06</u>	12.56	12.56
Handan	26.71	21.76	18.02	<u>15.01</u>	20.76	15.60	<b>14.95</b>
Hengshui	36.56	24.90	19.20	16.30	19.27	<b>15.20</b>	<u>16.07</u>
Langfang	21.45	12.39	14.23	<b>9.90</b>	13.21	10.21	<u>10.05</u>
Qinhuangdao	<u>17.18</u>	21.63	17.77	18.01	17.26	17.88	<b>15.58</b>
Shijiazhuang	25.85	20.44	22.12	20.53	21.49	<b>19.03</b>	<u>19.94</u>
Tangshan	35.01	29.71	<u>20.55</u>	22.10	21.09	<b>20.01</b>	20.99
Tianjin	18.29	11.38	13.90	<b>10.69</b>	14.51	<u>10.89</u>	11.26
Xingtai	31.02	22.43	19.88	<b>15.52</b>	21.56	16.28	<u>16.05</u>
Zhangjiakou	18.23	18.42	14.41	20.48	<u>13.29</u>	18.25	<b>12.62</b>
Average	25.26	18.89	17.43	15.69	17.15	<u>15.45</u>	<b>15.07</b>
City	FAC2						
	CMAQ	UK	RF	SVC	STVC	STAR	HDCM
Baoding	0.714	0.937	0.909	<b>0.960</b>	0.951	<b>0.960</b>	0.951
Beijing	0.644	<u>0.955</u>	0.824	0.953	0.918	<b>0.956</b>	0.944
Cangzhou	0.500	0.931	0.917	<b>0.967</b>	0.960	<b>0.967</b>	0.938
Chengde	0.630	0.829	0.851	<b>0.891</b>	0.777	0.832	<u>0.861</u>
Handan	0.823	0.889	0.957	0.981	0.910	0.984	<b>0.986</b>
Hengshui	0.536	0.862	0.953	<u>0.964</u>	0.924	<b>0.967</b>	0.949
Langfang	0.788	0.959	0.916	<b>0.992</b>	0.957	<b>0.992</b>	<b>0.992</b>
Qinhuangdao	0.674	0.624	0.694	0.715	0.720	<u>0.726</u>	<b>0.744</b>
Shijiazhuang	0.788	0.897	0.868	<u>0.911</u>	0.866	<b>0.918</b>	0.893
Tangshan	0.726	0.736	0.846	<u>0.781</u>	0.802	<b>0.879</b>	<u>0.851</u>
Tianjin	0.827	0.970	0.938	<b>0.981</b>	0.940	<u>0.980</u>	0.974
Xingtai	0.739	0.924	0.957	<u>0.995</u>	0.929	<b>0.997</b>	0.984
Zhangjiakou	0.636	0.679	<u>0.802</u>	0.674	0.750	0.677	<b>0.815</b>
Average	0.694	0.861	0.879	0.905	0.877	<u>0.910</u>	<b>0.914</b>

**S6.2. The performance of the proposed HDCM with spatial partitioning.** In this section, to illustrate that spatial partitioning works for the HDCM, we performed a cross-validation for the BTH data in both seasons of 2015 by considering the space-partitioning-based HDCM with two subregions (HDCM<sub>2</sub>).

TABLE S6

*Averaged RMSE (Root Mean Squared Error), CRPS (Continuous Rank Probability Score), MAE (Mean Absolute Error), and FAC2 (Fraction of Predictions within a factor of 2) for PM<sub>2.5</sub> concentration predictions ( $\mu\text{g}/\text{m}^3$ ) calculated for 13 cities in the BTH region using leave-one-city-out cross-validation from the the space-partitioning-based HDCM with two subregions (HDCM<sub>2</sub>). Daily data from the winter and summer of 2015 are considered.*

City	Winter of 2015				Summer of 2015			
	RMSE	CRPS	MAE	FAC2	RMSE	CRPS	MAE	FAC2
Baoding	70.03	37.99	48.44	0.90	22.58	11.92	16.34	0.95
Beijing	58.71	30.84	31.78	0.89	22.72	12.14	15.27	0.94
Cangzhou	41.05	22.34	26.27	0.95	18.34	10.08	13.85	0.94
Chengde	29.58	15.37	19.86	0.79	19.58	9.88	12.95	0.85
Handan	44.33	29.11	30.93	0.94	19.23	11.08	15.01	0.99
Hengshui	49.14	31.65	33.31	0.99	20.44	11.50	15.78	0.96
Langfang	37.67	29.03	23.61	0.99	13.83	9.63	9.96	0.99
Qinhuangdao	27.92	14.72	19.60	0.77	21.60	12.09	16.74	0.73
Shijiazhuang	53.89	33.30	38.77	0.94	26.16	14.44	20.28	0.90
Tangshan	38.28	22.48	27.03	0.95	25.88	14.94	21.28	0.87
Tianjin	36.70	26.41	23.81	0.96	14.84	8.94	11.35	0.97
Xingtai	47.46	32.48	32.46	0.96	21.36	12.01	16.01	0.98
Zhangjiakou	27.64	14.69	19.32	0.71	17.84	9.94	14.15	0.79
Average	43.26	26.19	28.86	0.90	20.34	11.43	15.31	0.91

**S7. A modified Diebold-Mariano test for predictive accuracy among methods.** To determine whether the difference between any two models is significant in terms of predictive accuracy, we carry out a modified Diebold-Mariano (DM) test (Harvey, Leybourne and Newbold, 1997).

Let  $\{r_t^{(1)}(\mathbf{s})\}_{t=1}^{N_t}$  and  $\{r_t^{(2)}(\mathbf{s})\}_{t=1}^{N_t}$  be the residual series at location  $\mathbf{s}$  obtained from Model 1 and Model 2 respectively. More specifically,  $r_t^{(1)}(\mathbf{s}) = \hat{y}_t^{(1)}(\mathbf{s}) - y_t(\mathbf{s})$  and  $r_t^{(2)}(\mathbf{s}) = \hat{y}_t^{(2)}(\mathbf{s}) - y_t(\mathbf{s})$ . Denote the loss-differential as  $d_{\mathbf{s},t} = g(r_t^{(1)}(\mathbf{s})) - g(r_t^{(2)}(\mathbf{s}))$ . Let  $\mu = E(d_{\mathbf{s},t})$  for a specific function  $g(\cdot)$ , where  $g(e) = e^2$  are used for this work. Under the null hypothesis  $\mu = 0$  (i.e., no difference in the accuracy of two sets of predictions), the DM test statistic (Diebold and Mariano, 1995) is given by

$$(S6) \quad \text{DM} = \frac{\bar{d}_s}{\sigma_{d_s}} \rightarrow \mathcal{N}(0, 1),$$

where  $\bar{d}_s = \frac{1}{N_t} \sum_{t=1}^{N_t} d_{\mathbf{s},t}$  and  $\sigma_{d_s}$  is the standard deviation of  $d_{\mathbf{s},t}$ .

In practice,  $\sigma_{d_s}$  is replaced by its consistent estimate. More specifically, a consistent estimate of  $\sigma_{d_s}$  can be asymptotically represented by  $\sqrt{\{\hat{\gamma}_s(0) - 2 \sum_{\tau=1}^{h-1} \hat{\gamma}_s(\tau)\}/N_t}$ , here  $\hat{\gamma}_s(\tau)$  is obtained by replacing  $\mu$  with  $\bar{d}_s$  into (S7):

$$(S7) \quad \gamma_s(\tau) = \frac{1}{N_t} \sum_{t=\tau+1}^{N_t} (d_{\mathbf{s},t} - \mu)(d_{\mathbf{s},t-\tau} - \mu).$$

where  $\tau = 0, 1, 2, \dots, N_t$ . (S7) defines the autocovariance at lag  $\tau$ .

Because the DM test (S6) tends to be oversized, a modified DM test is proposed by Harvey, Leybourne and Newbold (1997):

$$(S8) \quad \text{DM}^* = \sqrt{\frac{N_t + 1 - 2h + h(h-1)/N_t}{N_t}} \text{DM} \rightarrow \text{Student}(N_t - 1),$$

where  $h = N_t^{1/3} + 1$  is often used. In this work, we used the modified DM test ([S8](#)).

Because the dependence of the residuals of predictions is usually much weaker in space than in time, we assume the residual series,  $\{r_t^{(1)}(\mathbf{s})\}_{t=1}^{N_t}$  and  $\{r_t^{(2)}(\mathbf{s})\}_{t=1}^{N_t}$ , at different sites to be independent of each other. In light of this, the statistic  $DM^*$  used here is obtained by substituting the  $DM = \frac{\sum_{s=1}^{n_0} \bar{d}_s}{\sqrt{\sum_{s=1}^{n_0} \hat{\sigma}_{ds}^2}}$  into ([S8](#)), where  $n_0$  is the number of sites in the test set. Based on the predictions from LOOCV, the HLN statistics and their corresponding  $p$ -values are reported in Tables [S7](#) and [S8](#). In these tests, we always represent  $\{r_t^{(1)}(\mathbf{s})\}_{t=1}^{N_t}$  as the residuals of the proposed HDCM, and represent  $\{r_t^{(2)}(\mathbf{s})\}_{t=1}^{N_t}$  as that of one of the other seven methods.

Occasionally, even when HDCM is dominated by other methods in some cities under certain criteria (e.g., RMSE for the summer of 2015), the difference in the predictive accuracy is mostly not significant at a 5% significance level; see Table [S8](#) of the summer.

TABLE S7

The values of the modified Diebold-Mariano (DM) statistic and their corresponding p-values between the proposed HDCM and each of the other seven methods calculated for 13 cities in the BTH region using the results of leave-one-city-out cross-validation. A DM less than 0 indicates that the proposed HDCM performs better predictions than the method in the corresponding column, and the significant difference with the test level of 0.05 is represented by adding “\*” to the upper right corner of the statistic, and the corresponding p-value is in bold. A DM greater than 0 indicates that the method in the corresponding column performs better predictions than the proposed HDCM and is colored in gray, and the significant difference is represented by adding “†” to the upper right corner of the statistic. Daily data from November 1, 2015 to January 31, 2016 are considered.

City	Value of the modified DM statistic							
	CMAQ	UK	RF	SVC	STVC	STAR	HDCM	HDCM <sub>2</sub>
Baoding	-5.535*	-4.806*	-3.472*	-2.959*	-4.887*	-3.225*	-	5.151†
Beijing	-8.837*	-6.393*	-6.768*	-5.182*	-5.689*	-4.678*	-	-3.515*
Cangzhou	-4.623*	0.813	-0.494	0.982	1.016	1.568	-	-1.628
Chengde	-6.574*	-6.804*	-4.679*	-3.253*	-3.189*	-0.751	-	-3.491*
Handan	-8.199*	-5.736*	-3.262*	-4.060*	-3.405*	-4.985*	-	-0.182
Hengshui	-4.267*	-4.504*	-3.617*	-2.206*	-4.064*	-1.539	-	3.932†
Langfang	-7.963*	-2.973*	-6.706*	1.819	-3.313*	-0.945	-	-1.207
Qinhuangdao	-7.709*	-4.109*	-3.521*	-3.030*	-4.659*	-6.868*	-	-0.681
Shijiazhuang	-8.013*	2.449†	-0.174	3.080†	2.158†	3.485†	-	0.265
Tangshan	-9.542*	-4.135*	-5.992*	1.304	-3.105*	6.118†	-	-5.709*
Tianjin	-8.722*	-2.130*	-3.989*	2.142†	-2.253*	0.892	-	-4.568*
Xingtai	-4.873*	-3.388*	-3.901*	-2.315*	-3.232*	-2.029*	-	1.883
Zhangjiakou	-5.275*	-5.069*	-3.233*	-6.779*	-6.160*	-7.030*	-	2.146†

City	p-value of the modified DM test							
	CMAQ	UK	RF	SVC	STVC	STAR	HDCM	HDCM <sub>2</sub>
Baoding	<b>0.000</b>	<b>0.000</b>	<b>0.001</b>	<b>0.004</b>	<b>0.000</b>	<b>0.002</b>	-	0.000
Beijing	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	-	<b>0.001</b>
Cangzhou	<b>0.000</b>	0.418	0.623	0.329	0.312	0.120	-	0.107
Chengde	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.002</b>	<b>0.002</b>	0.455	-	<b>0.001</b>
Handan	<b>0.000</b>	<b>0.000</b>	<b>0.002</b>	<b>0.000</b>	<b>0.001</b>	<b>0.000</b>	-	0.856
Hengshui	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.030</b>	<b>0.000</b>	0.127	-	0.000
Langfang	<b>0.000</b>	<b>0.004</b>	<b>0.000</b>	0.072	<b>0.001</b>	0.347	-	0.231
Qinhuangdao	<b>0.000</b>	<b>0.000</b>	<b>0.001</b>	<b>0.003</b>	<b>0.000</b>	<b>0.000</b>	-	0.498
Shijiazhuang	<b>0.000</b>	0.016	0.862	0.003	0.034	0.001	-	0.791
Tangshan	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	0.196	<b>0.002</b>	0.000	-	<b>0.000</b>
Tianjin	<b>0.000</b>	<b>0.036</b>	<b>0.000</b>	0.035	<b>0.027</b>	0.375	-	<b>0.000</b>
Xingtai	<b>0.000</b>	<b>0.001</b>	<b>0.000</b>	<b>0.023</b>	<b>0.002</b>	<b>0.045</b>	-	0.063
Zhangjiakou	<b>0.000</b>	<b>0.000</b>	<b>0.002</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	-	0.035

The null hypothesis: Two predictions from two different models have the same accuracy.

TABLE S8

The values of the modified Diebold-Mariano (DM) statistic and their corresponding p-values between the proposed HDCM and each of the other seven methods calculated for 13 cities in the BTH region using the results of leave-one-city-out cross-validation. A DM less than 0 indicates that the proposed HDCM performs better predictions than the method in the corresponding column, and the significant difference with the test level of 0.05 is represented by adding “\*” to the upper right corner of the statistic, and the corresponding p-value is in bold. A DM greater than 0 indicates that the method in the corresponding column performs better predictions than the proposed HDCM and is colored in gray, and the significant difference is represented by adding “†” to the upper right corner of the statistic. Daily data from June 1, 2015, and August 31, 2015, are considered.

City	Value of the modified DM statistic							
	CMAQ	UK	RF	SVC	STVC	STAR	HDCM	HDCM <sub>2</sub>
Baoding	-10.868*	-2.847*	-3.425*	1.869	-1.450	-0.375	-	0.547
Beijing	<b>-12.224*</b>	2.005†	<b>-9.271*</b>	3.192†	<b>-5.222*</b>	<b>-1.503</b>	-	-1.497
Cangzhou	-8.760*	-3.143*	-1.685	3.892†	0.587	3.648†	-	2.534†
Chengde	-4.398*	0.062	-1.432	1.355	1.450	1.160	-	-2.687*
Handan	-8.460*	-8.725*	<b>-5.105*</b>	-0.038	<b>-9.187*</b>	<b>-1.628</b>	-	-0.591
Hengshui	-11.840*	-6.984*	-3.810*	0.022	-3.994*	1.457	-	2.765†
Langfang	-8.617*	-2.812*	<b>-5.907*</b>	1.736	<b>-4.460*</b>	0.539	-	-0.353
Qinhuangdao	-4.483*	-5.211*	-4.361*	<b>-3.910*</b>	<b>-3.817*</b>	<b>-2.587*</b>	-	-2.213*
Shijiazhuang	-5.745*	-1.320	<b>-3.457*</b>	0.511	<b>-1.990*</b>	1.159	-	0.391
Tangshan	-5.208*	-6.814*	0.427	<b>-2.970*</b>	0.133	3.890†	-	-1.078
Tianjin	-8.202*	-0.720	<b>-4.889*</b>	1.786	<b>-5.861*</b>	0.667	-	-1.237
Xingtai	-9.456*	-7.077*	<b>-6.951*</b>	1.264	<b>-4.469*</b>	<b>-1.397</b>	-	-1.519
Zhangjiakou	-5.733*	-7.285*	<b>-5.552*</b>	<b>-7.211*</b>	<b>-2.596*</b>	<b>-6.641*</b>	-	-5.156*

City	p-value of the modified DM test							
	CMAQ	UK	RF	SVC	STVC	STAR	HDCM	HDCM <sub>2</sub>
Baoding	<b>0.000</b>	<b>0.005</b>	<b>0.001</b>	0.065	0.151	0.709	-	0.586
Beijing	<b>0.000</b>	0.048	<b>0.000</b>	0.002	<b>0.000</b>	0.136	-	0.138
Cangzhou	<b>0.000</b>	<b>0.002</b>	0.095	0.000	0.559	0.000	-	0.013
Chengde	<b>0.000</b>	0.951	0.156	0.179	0.150	0.249	-	<b>0.009</b>
Handan	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	0.970	<b>0.000</b>	0.107	-	0.556
Hengshui	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	0.982	<b>0.000</b>	0.149	-	0.007
Langfang	<b>0.000</b>	<b>0.006</b>	<b>0.000</b>	0.086	<b>0.000</b>	0.592	-	0.725
Qinhuangdao	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.011</b>	-	<b>0.029</b>
Shijiazhuang	<b>0.000</b>	0.190	<b>0.001</b>	0.610	<b>0.050</b>	0.250	-	0.697
Tangshan	<b>0.000</b>	<b>0.000</b>	0.670	<b>0.004</b>	0.894	0.000	-	0.284
Tianjin	<b>0.000</b>	0.473	<b>0.000</b>	0.077	<b>0.000</b>	0.506	-	0.219
Xingtai	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	0.210	<b>0.000</b>	0.166	-	0.132
Zhangjiakou	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.000</b>	<b>0.011</b>	<b>0.000</b>	-	<b>0.000</b>

## S8. Robustness of the proposed approach to deviations from the normality and linearity assumptions.

**S8.1. Robustness to normality.** Robustness to normality is investigated through a cross validation for the raw data based on the proposed models (i.e., HDCM and HDCM<sub>2</sub>). For the BTH PM<sub>2.5</sub> concentration data, the raw observations in the summer of 2015 tend to admit a skewed right distribution where the right tail (larger values) is much longer than the left tail (smaller values), while the distribution of observations on the square root scale from the summer is very close to a normal distribution; see Figure S1(d) and Figure S1(e). Therefore, we perform a cross-validation for the raw observations in the summer to verify whether our method is valid when the normality assumption is violated. LOCO CV results are reported in Table S9.

TABLE S9

*Averaged RMSE (Root Mean Squared Error), CRPS (Continuous Rank Probability Score), MAE (Mean Absolute Error), and FAC2 (Fraction of Predictions within a factor of 2) for PM<sub>2.5</sub> concentration predictions ( $\mu\text{g}/\text{m}^3$ ) calculated for 13 cities in the BTH region using leave-one-city-out cross-validation from eight models: (1) CMAQ numerical model output, (2) universal kriging, (3) random forest model, (4) spatially-varying coefficient downscaling model (SVC), (5) spatiotemporally-varying coefficient downscaling model (STVC), (6) first-order spatiotemporal autoregression (STAR), (7) proposed hierarchical dynamic calibration model (HDCM), and (8) spatial-partitioning-based HDCM with two subregions (HDCM<sub>2</sub>). The smallest RMSE, CRPS, MAE, and FAC2 are in bold and the second smallest ones are underlined. Daily raw data from June 1, 2015, and August 31, 2015, are used for all models.*

Criterion	CMAQ	UK	RF	SVC	STVC	STAR	HDCM	HDCM <sub>2</sub>
RMSE	32.97	24.67	23.09	22.03	22.79	20.46	<b>20.35</b>	20.43
CRPS	19.54	13.95	<u>12.57</u>	12.15	12.70	11.33	<u>11.31</u>	<b>11.25</b>
MAE	25.26	19.12	17.69	16.75	17.58	15.52	<u>15.51</u>	<b>15.41</b>
FAC2	0.64	0.86	0.87	0.89	0.85	<u>0.90</u>	0.89	<b>0.90</b>

**S8.2. Robustness to the linearity.** The PM<sub>2.5</sub> concentrations are often nonlinearly affected by different weather conditions such as temperature, wind power (Liang et al., 2015). We investigate the robustness of the proposed approach to the linearity by comparing the proposed models with the additive model (ADM). We use ADM of the form

$$(S9) \quad y_t(\mathbf{s}) = \beta_0 + \beta_1 x_{t,1}(C_s) + \sum_{\kappa=2}^4 g_{\kappa-1}(x_{t,\kappa}(C_s)) + g_4(s^x, s^y, t) + \varepsilon_t(s),$$

where  $x_{t,1}(C_s)$  is the weighted outputs of CMAQ at location  $s$  and at time point  $t$ ,  $\{x_{t,2}(C_s), x_{t,3}(C_s), x_{t,4}(C_s)\}$  represent the weighted surface temperature, northern cumulative wind power, and eastern cumulative wind power, respectively; and  $s^x$  and  $s^y$  denote longitude and latitude, respectively. The functions  $\{g_k(\cdot)\}, k = 1, \dots, 4$  are estimated by basis expansion  $g_k(z_{t,k}(s)) = \sum_{j=1}^{J_k} \alpha_{k,j} \phi_{k,j}$  with  $J_k = 5$  for univariate functions  $\{g_k(\cdot)\}_{k=1}^3$  and  $J_4 = 180$  for the three-dimensional smooth function  $g_4(\cdot)$ . We use cyclic cubic-spline basis functions for  $\{\phi_{k,j}\}_{k=1}^3$  and thin plate splines for  $\phi_{4,j}$ . See Wood (2017) for more discussions on additive models. ADM is implemented using the R package mgcv (Wood, 2022), and the results of cross-validation with ADM (S9) are presented in Table S10.

ADM (S9) is used to fit the BTH PM<sub>2.5</sub> concentration data in the summer of 2015. Although relationships between the PM<sub>2.5</sub> readings and other meteorological variables are nonlinear, the average value of the different metrics across all cities shows that HDCM (or HDCM<sub>2</sub>) using the linear trends outperforms ADM (S9) using the nonlinear trends; see Tables S4-S6 for the performance of HDCM and Table S10 for ADM. Therefore, the proposed approach is robust to the deviation from the linearity assumption on fixed effects.

TABLE S10

*Averaged RMSE (Root Mean Squared Error), CRPS (Continuous Rank Probability Score), MAE (Mean Absolute Error), and FAC2 (Fraction of Predictions within a factor of 2) for PM<sub>2.5</sub> concentration predictions ( $\mu\text{g}/\text{m}^3$ ) calculated for 13 cities in the BTH region using leave-one-city-out cross-validation from the additive model (ADM). Daily data from November 1, 2015 to January 31, 2016 are considered. Daily data from June 1, 2015, and August 31, 2015, are considered.*

City	RMSE	CRPS	MAE	FAC2
Baoding	24.24	12.97	17.60	0.95
Beijing	24.86	13.30	17.55	0.90
Cangzhou	18.69	10.62	14.65	0.94
Chengde	18.80	10.32	14.05	0.77
Handan	24.40	13.81	19.54	0.96
Hengshui	23.19	12.51	17.18	0.96
Langfang	17.85	10.25	13.74	0.95
Qinhuangdao	25.90	13.66	17.45	0.66
Shijiazhuang	27.47	15.33	21.90	0.87
Tangshan	25.39	14.60	21.04	0.89
Tianjin	19.03	10.92	15.35	0.90
Xingtai	24.13	13.25	18.22	0.97
Zhangjiakou	43.18	24.14	34.90	0.47
Average	24.39	13.51	18.70	0.86

## S9. Variational Bayes for HDCM parameters.

S9.1. *Some additional results from spEnKS.* By Algorithm 1 of the Manuscript, we have  $q_{v_t}(\mathbf{v}_t) \sim \mathcal{N}\left(\widehat{\mathbf{v}}_{t|N_t}, \widehat{\boldsymbol{\Sigma}}_{t,t|N_t}\right)$  for  $t \leq N_t$ . Furthermore, the following results are necessary to run the VB procedure from (S15) to (S22):

$$(S10) \quad \mathbf{S}_{00}^{(r)} = \sum_{t=1}^{N_t} \langle \mathbf{v}_{t-1}^{(r)} \mathbf{v}_{t-1}^{T,(r)} \rangle = \sum_{t=1}^{N_t} \left( \widehat{\mathbf{v}}_{t-1|N_t}^{(r)} \widehat{\mathbf{v}}_{t-1|N_t}^{T,(r)} + \widehat{\boldsymbol{\Sigma}}_{t-1,t-1|N_t}^{(r)} \right),$$

$$(S11) \quad \mathbf{S}_{10}^{(r)} = \sum_{t=1}^{N_t} \langle \mathbf{v}_t^{(r)} \mathbf{v}_{t-1}^{T,(r)} \rangle = \sum_{t=1}^{N_t} \left( \widehat{\mathbf{v}}_{t|N_t}^{(r)} \widehat{\mathbf{v}}_{t-1|N_t}^{T,(r)} + \widehat{\boldsymbol{\Sigma}}_{t,t-1|N_t}^{(r)} \right),$$

$$(S12) \quad \mathbf{S}_{11}^{(r)} = \sum_{t=1}^{N_t} \langle \mathbf{v}_t^{(r)} \mathbf{v}_t^{T,(r)} \rangle = \sum_{t=1}^{N_t} \left( \widehat{\mathbf{v}}_{t|N_t}^{(r)} \widehat{\mathbf{v}}_{t|N_t}^{T,(r)} + \widehat{\boldsymbol{\Sigma}}_{t,t|N_t}^{(r)} \right),$$

$$(S13) \quad \mathbf{S}_{11} = \text{diag}\left(\mathbf{S}_{11}^{(1)}, \dots, \mathbf{S}_{11}^{(R)}\right).$$

S9.2. *Coordinate ascent mean-field VB for HDCM.* Following Blei, Kucukelbir and McAuliffe (2017), we decompose  $\log p(\mathbf{y})$  into the following two terms:

$$\begin{aligned} \log p(\mathbf{y}) &= \log p(\mathbf{y}, \boldsymbol{\Theta}, \mathbf{v}) - \log p(\boldsymbol{\Theta}, \mathbf{v}|\mathbf{y}) \\ &= \log \left\{ \frac{p(\mathbf{y}, \boldsymbol{\Theta}, \mathbf{v})}{q(\boldsymbol{\Theta}, \mathbf{v})} \right\} - \log \left\{ \frac{p(\boldsymbol{\Theta}, \mathbf{v}|\mathbf{y})}{q(\boldsymbol{\Theta}, \mathbf{v})} \right\} \\ &= \int q(\boldsymbol{\Theta}, \mathbf{v}) \log \left\{ \frac{p(\mathbf{y}, \boldsymbol{\Theta}, \mathbf{v})}{q(\boldsymbol{\Theta}, \mathbf{v})} \right\} - \int q(\boldsymbol{\Theta}, \mathbf{v}) \log \left\{ \frac{p(\boldsymbol{\Theta}, \mathbf{v}|\mathbf{y})}{q(\boldsymbol{\Theta}, \mathbf{v})} \right\} \\ (S14) \quad &= \mathcal{L}(q) + \text{KL}[q(\boldsymbol{\Theta}, \mathbf{v}) \parallel p(\boldsymbol{\Theta}, \mathbf{v}|\mathbf{y})] \geq \mathcal{L}(q). \end{aligned}$$

We approximate the posterior using  $q \equiv q(\Theta, v)$ . The objective of optimization is to minimize the Kullback-Leibler divergence  $\text{KL}[q(\Theta, v) | p(\Theta, v | \mathbf{y})]$ , or equivalently, to maximize the function  $\mathcal{L}(q)$  with respect to the function  $q$ .

Using the mean-field approximation method (Blei, Kucukelbir and McAuliffe, 2017), we assume  $q(\Theta, v)$  in (S14) is separable and obtain  $q = q_\Theta(\Theta)q_v(v)$  where  $q_\Theta(\Theta) = q_\beta(\beta)q_{\sigma^2}(\sigma^2)q_{\tau_0^2}(\tau_0^2)q_{\tau^2}(\tau^2)q_{\zeta_0^2}(\zeta_0^2)q_{\zeta^2}(\zeta^2)q_{\theta_1}(\theta_1)q_{\theta_2}(\theta_2)$ . Note that  $q_v(v)$  is approximated using the smoothing ensemble from spEnKS of Section 3.2.4 of the Manuscript, and a complete coordinate ascent procedure for inferring  $q_\Theta(\Theta)$  is described as follows:

- a. For the regression coefficient  $\beta$ , we set  $\pi(\beta) = \mathcal{N}(\mu_\beta, \Sigma_\beta)$  and have

$$(S15) \quad q_\beta(\beta) = \mathcal{N}(\mathbf{D}_\beta \boldsymbol{\xi}_\beta, \mathbf{D}_\beta),$$

where  $\mathbf{D}_\beta = \left( \langle \frac{1}{\sigma^2} \rangle \sum_{t=1}^{N_t} \mathbf{X}_t^T \mathbf{X}_t + \Sigma_\beta^{-1} \right)^{-1}$  and  $\boldsymbol{\xi}_\beta = \langle \frac{1}{\sigma^2} \rangle \sum_{t=1}^{N_t} \mathbf{X}_t^T (\mathbf{y}_t - \mathbf{H}\hat{\mathbf{v}}_{t|N_t}) + \Sigma_\beta^{-1} \mu_\beta$ .

- b. For variance of the noise  $\sigma^2$ , we set  $\pi(\sigma^2) = IG(a_{\sigma^2}, b_{\sigma^2})$  with the shape  $a_{\sigma^2}$  and the scale  $b_{\sigma^2}$ , then

$$(S16) \quad q_{\sigma^2}(\sigma^2) = IG\left(a_{\sigma^2} + nN_t/2, b_{\sigma^2} + f_{\sigma^2}/2\right),$$

where ‘‘IG’’ denotes an Inverse Gamma distribution,  $f_{\sigma^2} = \sum_{t=1}^{N_t} \left[ \mathbf{y}_t^T (\mathbf{y}_t - 2\mathbf{X}_t \langle \beta \rangle) - 2\hat{\mathbf{v}}_t^T \mathbf{H}^T (\mathbf{y}_t - \mathbf{X}_t \langle \beta \rangle) + \text{Tr}\left\{ \mathbf{X}_t^T \mathbf{X}_t (\mathbf{D}_\beta + \langle \beta \beta^T \rangle) \right\} \right] + \text{Tr}(\mathbf{H}^T \mathbf{H} \mathbf{S}_{11})$ , and  $\mathbf{S}_{11}$  is defined in (S13), here  $\text{Tr}(\mathbf{A})$  denotes the trace of the matrix  $\mathbf{A}$ .

- c. For the process transition parameter  $\theta_1^{(r)}$  of IDE (4) of the Manuscript at the  $r$ th subregion, we set  $\pi(\theta_1^{(r)}) = \mathcal{N}(\mu_{\theta_1^{(r)}}, \sigma_{\theta_1^{(r)}}^2)$ , then

$$(S17) \quad q_{\theta_1^{(r)}}(\theta_1^{(r)}) = \mathcal{N}\left(D_{\theta_1^{(r)}} \boldsymbol{\xi}_{\theta_1^{(r)}}, D_{\theta_1^{(r)}}\right),$$

where  $D_{\theta_1} = \left\{ \text{Tr}\left(\mathbf{M}_{\langle \theta_2^{(r)} \rangle} \langle \mathbf{Q}^{(r)} \rangle \mathbf{M}_{\langle \theta_2^{(r)} \rangle} \mathbf{S}_{00}^{(r)}\right) + \sigma_{\theta_1^{(r)}}^{-2} \right\}^{-1}$ ,  $\xi_{\theta_1} = \text{Tr}\left(\mathbf{M}_{\langle \theta_2^{(r)} \rangle} \langle \mathbf{Q}^{(r)} \rangle \mathbf{S}_{10}^{(r)}\right) + \mu_{\theta_1^{(r)}} \sigma_{\theta_1^{(r)}}^{-2}$ ,  $\mathbf{S}_{00}^{(r)}$  and  $\mathbf{S}_{10}^{(r)}$  are defined in (S10) and (S11), respectively.

- d. For the scale parameter  $\tau^{2,(r)}$  of IDE (4) of the Manuscript at the  $r$ th subregion, let  $\pi(\tau^{2,(r)}) = \text{Gamma}(a_{\tau^{2,(r)}}, b_{\tau^{2,(r)}})$  with the shape  $a_{\tau^{2,(r)}}$  and the scale  $b_{\tau^{2,(r)}}$ . We have

$$(S18) \quad q_{\tau^{2,(r)}}(\tau^{2,(r)}) = \text{Gamma}(\tilde{a}_{\tau^{2,(r)}}, \tilde{b}_{\tau^{2,(r)}}),$$

with  $\tilde{a}_{\tau^{2,(r)}} = a_{\tau^{2,(r)}} + m_r N_t / 2$ ,  $\tilde{b}_{\tau^{2,(r)}} = b_{\tau^{2,(r)}} + f_{\tau^{2,(r)}} / 2$ , where let  $\langle \Lambda^{(r)} \rangle = \mathbf{G}^{(r)} + \langle \zeta^{2,(r)} \rangle \mathbf{I}_{m_r}$  and

$$f_{\tau^{2,(r)}} = \text{Tr}\left(\langle \Lambda^{(r)} \rangle \mathbf{S}_{11}^{(r)} - 2\langle \tilde{\theta}_1^{(r)} \rangle \mathbf{M}_{\langle \theta_2^{(r)} \rangle} \langle \Lambda^{(r)} \rangle \mathbf{S}_{10}^{(r)} + \langle \tilde{\theta}_1^{2,(r)} \rangle \mathbf{M}_{\langle \theta_2^{(r)} \rangle} \langle \Lambda^{(r)} \rangle \mathbf{M}_{\langle \theta_2^{(r)} \rangle} \mathbf{S}_{00}^{(r)}\right).$$

- e. For the scale parameter  $\tau_0^{2,(r)}$  of the initial state of IDE (4) of the Manuscript at the  $r$ th subregion, we set  $\pi(\tau_0^{2,(r)}) = \text{Gamma}(a_{\tau_0^{2,(r)}}, b_{\tau_0^{2,(r)}})$  and have

$$(S19) \quad q_{\tau_0^{2,(r)}}(\tau_0^{2,(r)}) = \text{Gamma}(\tilde{a}_{\tau_0^{2,(r)}}, \tilde{b}_{\tau_0^{2,(r)}}),$$

with  $\tilde{a}_{\tau_0^{2,(r)}} = a_{\tau_0^{2,(r)}} + m_r / 2$ ,  $\tilde{b}_{\tau_0^{2,(r)}} = b_{\tau_0^{2,(r)}} + \text{Tr}\{\langle \Lambda_0^{(r)} \rangle \langle \mathbf{v}_0 \mathbf{v}_0^T \rangle\} / 2$ , here  $\langle \Lambda_0^{(r)} \rangle = \mathbf{G}^{(r)} + \langle \zeta_0^{2,(r)} \rangle \mathbf{I}_{m_r}$ .

- f. Although  $q_{\theta_2^{(r)}}$ ,  $q_{\zeta^{2,(r)}}$ , and  $q_{\zeta_0^{2,(r)}}$  do not have closed form expressions, they are proportional to:

$$(S20) \quad q_{\theta_2}(\theta_2^{(r)}) \propto \exp(-f_{\theta_2^{(r)}}/2)\pi(\theta_2^{(r)}),$$

$$(S21) \quad q_{\zeta^{2,(r)}}(\zeta^{2,(r)}) \propto |\mathbf{Q}^{(r)}|^{\frac{N_t}{2}} \exp(-f_{\zeta^{2,(r)}}/2)\pi(\zeta^{(r)}),$$

$$(S22) \quad q_{\zeta_0^{2,(r)}}(\zeta_0^{2,(r)}) \propto |\mathbf{Q}_0^{(r)}|^{\frac{1}{2}} \exp(-f_{\zeta_0^{2,(r)}}/2)\pi(\zeta_0^{(r)}),$$

where

$$\begin{aligned} f_{\theta_2} &= \text{Tr} \left\{ -2\langle \theta_1^{(r)} \rangle \mathcal{M}_{\theta_2^{(r)}} \langle \mathbf{Q}^{(r)} \rangle \mathbf{S}_{10}^{(r)} + (\langle \theta_1^{2,(r)} \rangle + D_{\theta_1^{(r)}}) \mathcal{M}_{\theta_2^{(r)}} \langle \mathbf{Q}^{(r)} \rangle \mathcal{M}_{\theta_2^{(r)}} \mathbf{S}_{00}^{(r)} \right\}, \\ f_{\zeta^{2,(r)}} &= \text{Tr} \left\{ \mathbf{S}_{11}^{(r)} - 2\langle \theta_1^{(r)} \rangle \mathcal{M}_{\langle \theta_2^{(r)} \rangle} \mathbf{S}_{10}^{(r)} + (\langle \theta_1^{2,(r)} \rangle + D_{\theta_1^{(r)}}) \mathcal{M}_{\langle \theta_2^{(r)} \rangle} \mathcal{M}_{\langle \theta_2^{(r)} \rangle} \mathbf{S}_{00}^{(r)} \right\} \langle \tau^{2,(r)} \rangle \zeta^{2,(r)}, \\ f_{\zeta_0^{2,(r)}} &= \text{Tr} \left\{ (\hat{\mathbf{v}}_0^{(r)} \hat{\mathbf{v}}_0^{T,(r)} + \hat{\Sigma}_{0,0|N_t}^{(r)}) \right\} \langle \tau_0^{2,(r)} \rangle \zeta_0^{2,(r)}. \end{aligned}$$

Based on the results above, these three  $q$  functions can be approximated with the Laplace method (Section 4.7 of [Lindgren and Rue \(2015\)](#)). Section S10 provides further details on the Laplace approximation related to the proposed approach.

### S10. The Laplace approximation for the parameter without the conjugate prior.

Assume the function  $f(\xi)$  is the density of a random variable  $\xi$ . Based on the Laplace method, the integral of  $f(\xi)$  can be approximated by a Gaussian integral. Specifically, the mean of the Gaussian distribution is a stationary point of  $\log f(\xi)$  such that  $\frac{\partial \log f(\xi)}{\partial \xi} = 0$ , denoted by  $\xi_0$ , and the variance is  $-1/\frac{\partial^2 \log f(\xi)}{\partial^2 \xi}|_{\xi=\xi_0}$ , where the vertical bar indicates a specific value of a parameter. See Section 4.7 of [Lindgren and Rue \(2015\)](#) for more details on this topic.

Several parameters in the IDE (4) of the Manuscript do not have the conditionally conjugate prior, including  $\theta_2$ ,  $\zeta_0$ , and  $\zeta$ , leading to no closed-form expressions for their  $q$  functions. But because these functions are proportional to (S20)-(S22) of the manuscript, we can approximate their posterior distributions with the Gaussian by finding corresponding maximum posterior points by (S20), (S21), and (S22), respectively. Below are the relevant details:

- Find the maximum posterior estimate  $\theta_2^{*,(r)}$  for  $\theta_2^{(r)}$  by solving  $\frac{\partial \log \pi(\theta_2^{(r)})}{\partial \theta_2^{(r)}} - \frac{1}{2} \frac{\partial f_{\theta_2^{(r)}}}{\partial \theta_2^{(r)}} = 0$ , where  $\frac{\partial \log \pi(\theta_2^{(r)})}{\partial \theta_2^{(r)}} = 0$  for an uniform distribution  $\pi(\theta_2^{(r)})$ . The second term on the left-hand side of the equation above can be written as

$$\frac{\partial f_{\theta_2^{(r)}}}{\partial \theta_2^{(r)}} = \text{Tr} \left\{ -2\langle \theta_1^{(r)} \rangle \frac{\partial \mathcal{M}_{\theta_2^{(r)}}}{\partial \theta_2^{(r)}} \langle \mathbf{Q}^{(r)} \rangle \mathbf{S}_{10}^{(r)} + (\langle \theta_1^{2,(r)} \rangle + D_{\theta_1^{(r)}}) \frac{\partial f_{\theta_2^{(r)}}^{\mathcal{M}}}{\partial \theta_2^{(r)}} \mathbf{S}_{00}^{(r)} \right\},$$

where the  $(l, l')$ th element of  $\frac{\partial \mathcal{M}_{\theta_2^{(r)}}}{\partial \theta_2^{(r)}}$  is  $\frac{\partial W(d_{l,l'}; \theta_2^{(r)})}{\partial \theta_2^{(r)}} = \frac{d_{l,l'}}{4\theta_2^{2,(r)}} \{ (1 - d_{l,l'}/\theta_2^{(r)})^2 (1 + 3d_{l,l'}/\theta_2^{(r)}) - (1 - d_{l,l'}/\theta_2^{(r)})^3 \}$ ,  $f_{\theta_2^{(r)}}^{\mathcal{M}} = \mathcal{M}_{\theta_2^{(r)}} \langle \mathbf{Q}^{(r)} \rangle \mathcal{M}_{\theta_2^{(r)}}$ , and the derivative of  $f_{\theta_2^{(r)}}^{\mathcal{M}}$  is given by  $\frac{\partial f_{\theta_2^{(r)}}^{\mathcal{M}}}{\partial \theta_2^{(r)}} = \frac{\partial \mathcal{M}_{\theta_2^{(r)}}}{\partial \theta_2^{(r)}} \langle \mathbf{Q}^{(r)} \rangle \mathcal{M}_{\theta_2^{(r)}} + \mathcal{M}_{\theta_2^{(r)}} \langle \mathbf{Q}^{(r)} \rangle \frac{\partial \mathcal{M}_{\theta_2^{(r)}}}{\partial \theta_2^{(r)}}$ . Moreover,  $\frac{\partial^2 f_{\theta_2^{(r)}}}{\partial \theta_2^{(r)}} = \text{Tr} \left\{ -2\langle \theta_1 \rangle \frac{\partial^2 \mathcal{M}_{\theta_2^{(r)}}}{\partial \theta_2^{(r)}} \langle \mathbf{Q}^{(r)} \rangle \mathbf{S}_{10}^{(r)} + (\langle \theta_1^{2,(r)} \rangle + D_{\theta_1^{(r)}}) \frac{\partial^2 f_{\theta_2^{(r)}}^{\mathcal{M}}}{\partial \theta_2^{(r)}} \mathbf{S}_{00}^{(r)} \right\}$ , where  $\frac{\partial^2 f_{\theta_2^{(r)}}^{\mathcal{M}}}{\partial \theta_2^{(r)}} = \frac{\partial^2 \mathcal{M}_{\theta_2^{(r)}}}{\partial \theta_2^{(r)}} \langle \mathbf{Q}^{(r)} \rangle \mathcal{M}_{\theta_2^{(r)}} + 2 \frac{\partial \mathcal{M}_{\theta_2^{(r)}}}{\partial \theta_2^{(r)}} \langle \mathbf{Q}^{(r)} \rangle \frac{\partial \mathcal{M}_{\theta_2^{(r)}}}{\partial \theta_2^{(r)}} + \mathcal{M}_{\theta_2^{(r)}} \langle \mathbf{Q}^{(r)} \rangle \frac{\partial^2 \mathcal{M}_{\theta_2^{(r)}}}{\partial \theta_2^{(r)}}$ , here the  $(l, l')$ th element of  $\frac{\partial^2 \mathcal{M}_{\theta_2^{(r)}}}{\partial \theta_2^{(r)}}$  is  $\frac{\partial^2 W(d_{l,l'}; \theta_2^{(r)})}{\partial \theta_2^{(r)}} = -\frac{d_{l,l'}}{2\theta_2^{3,(r)}} \{ (1 - d_{l,l'}/\theta_2^{(r)})^2 (1 + 3d_{l,l'}/\theta_2^{(r)}) - (1 - d_{l,l'}/\theta_2^{(r)})^3 \} +$

$\frac{d_{l,l'}^2}{2\theta_2^{4,(r)}} \left\{ (1 - d_{l,l'} / \theta_2^{(r)})^2 (1 + 3d_{l,l'} / \theta_2^{(r)}) - 3(1 - d_{l,l'} / \theta_2^{(r)})^2 \right\}$ . We then have  $\theta_2^{(r)} \sim \mathcal{N}(\theta_2^{*,(r)}, -1 / \left. \frac{\partial^2 f_{\theta_2^{(r)}}}{\partial \theta_2^{(r)}} \right|_{\theta_2^{(r)} = \theta_2^{*,(r)}})$ .

2. Find the maximum posteriori estimate  $\zeta^{2*,(r)}$  for  $\zeta^{2,(r)}$  by solving  $N_t \frac{\partial \log |\mathbf{G}^{(r)} + \zeta^{2,(r)} \mathbf{I}_{m_r}|}{\partial \zeta^{2,(r)}} - \frac{\partial f_{\zeta^{2,(r)}}}{\partial \zeta^{2,(r)}} + \frac{\partial \log \pi(\zeta^{2,(r)})}{\partial \zeta^{2,(r)}} = 0$ , where  $\frac{\partial \log |\mathbf{G}^{(r)} + \zeta^{2,(r)} \mathbf{I}_{m_r}|}{\partial \zeta^{2,(r)}} = \text{Tr}\left\{(\mathbf{G}^{(r)} + \zeta^{2,(r)} \mathbf{I}_{m_r})^{-1}\right\}$ ,  $\frac{\partial f_{\zeta^{2,(r)}}}{\partial \zeta^{2,(r)}} = \text{Tr}\left\{\mathbf{S}_{11}^{(r)} - 2\langle \theta_1^{(r)} \rangle \mathbf{M}_{\langle \theta_2^{(r)} \rangle} \mathbf{S}_{10}^{(r)} + (\langle \theta_1^{2,(r)} \rangle + D_{\theta_1^{(r)}}) \mathbf{M}_{\langle \theta_2^{(r)} \rangle} \mathbf{M}_{\langle \theta_2^{(r)} \rangle} \mathbf{S}_{00}^{(r)}\right\} \langle \tau^{2,(r)} \rangle$ , and  $\frac{\partial \log \pi(\zeta^{2,(r)})}{\partial \zeta^{2,(r)}} = 0$ . On the other hand, we find  $\sigma_{\zeta^{2*,(r)}}^2 = -1 / \left. \frac{\partial^2 \log |\mathbf{G}^{(r)} + \zeta^{2,(r)} \mathbf{I}_{m_r}|}{\partial \zeta^{4,(r)}} \right|_{\sigma_{\zeta^{2,(r)}}^2 = \zeta^{2*,(r)}}$ , where  $\frac{\partial^2 \log |\mathbf{G}^{(r)} + \zeta^{2,(r)} \mathbf{I}_{m_r}|}{\partial \zeta^{4,(r)}} = -(\mathbf{G}^{(r)} + \zeta^{2,(r)} \mathbf{I}_{m_r})^{-1} (\mathbf{G}^{(r)} + \zeta^{2,(r)} \mathbf{I}_{m_r})^{-1}$  by using  $\frac{\partial \text{Tr} \tilde{\mathbf{G}}^{(r)}}{\partial \tilde{\mathbf{G}}^{(r)}} = \mathbf{I}_{m_r}$  and  $\frac{\partial \tilde{\mathbf{G}}^{-1,(r)}}{\partial \zeta^{2,(r)}} = -\tilde{\mathbf{G}}^{-1,(r)} \frac{\partial \tilde{\mathbf{G}}^{(r)}}{\partial \zeta^{2,(r)}} \tilde{\mathbf{G}}^{-1,(r)}$ , here  $\tilde{\mathbf{G}}^{(r)} = \mathbf{G}^{(r)} + \zeta^{2,(r)} \mathbf{I}_{m_r}$ .
3. Find the maximum posteriori estimate  $\zeta_0^{2*,(r)}$  for  $\zeta_0^{2,(r)}$  by solving  $\frac{\partial \log |\mathbf{G}^{(r)} + \zeta_0^{2,(r)} \mathbf{I}_{m_r}|}{\partial \zeta_0^{2,(r)}} - \frac{\partial f_{\zeta_0^{2,(r)}}}{\partial \zeta_0^{2,(r)}} + \frac{\partial \log \pi(\zeta_0^{2,(r)})}{\partial \zeta_0^{2,(r)}} = 0$ , where  $\frac{\partial \log |\mathbf{G}^{(r)} + \zeta_0^{2,(r)} \mathbf{I}_{m_r}|}{\partial \zeta_0^{2,(r)}} = \text{Tr}\left\{(\mathbf{G}^{(r)} + \zeta_0^{2,(r)} \mathbf{I}_{m_r})^{-1}\right\}$ ,  $\frac{\partial f_{\zeta_0^{2,(r)}}}{\partial \zeta_0^{2,(r)}} = \text{Tr}\left\{(\hat{\mathbf{v}}_0 \hat{\mathbf{v}}_0^T + \hat{\Sigma}_{0,0|N_t})\right\} \langle \tau_0^{2,(r)} \rangle$ , and  $\frac{\partial \log \pi(\zeta_0^{2,(r)})}{\partial \zeta_0^{2,(r)}} = 0$ . On the other hand, we find  $\sigma_{\zeta_0^{2*,(r)}}^2 = -1 / \left. \frac{\partial^2 \log |\mathbf{G}^{(r)} + \zeta_0^{2,(r)} \mathbf{I}_{m_r}|}{\partial \zeta_0^{4,(r)}} \right|_{\sigma_{\zeta_0^{2,(r)}}^2 = \zeta_0^{2*,(r)}}$ .

In practice, when optimization is used to find a stationary point, to make the algorithm more stable, the logarithmic transformation may be useful for these three parameters, i.e.,  $\log \theta_2^{(2)} = \vartheta_1$ ,  $\log \zeta_0^{2,(r)} = \vartheta_2$ , and  $\log \zeta^{2,(r)} = \vartheta_3$ , while other details remain the same as above.

## REFERENCES

- BAKAR, K. S., KOKIC, P. and JIN, H. (2016). Hierarchical spatially varying coefficient and temporal dynamic process models using spTDyn. *J. Stat. Comput. Simul.* **86** 820–840.
- BERROCAL, V. J., GELFAND, A. E. and HOLLAND, D. M. (2010). A spatio-temporal downscaler for output from numerical models. *J. Agric. Biol. Environ. Stat.* **15** 176–197.
- BERROCAL, V. J., GUAN, Y., MUYSKENS, A., WANG, H., REICH, B. J., MULHOLLAND, J. A. and CHANG, H. H. (2020). A comparison of statistical and machine learning methods for creating national daily maps of ambient PM<sub>2.5</sub> concentration. *Atmos. Environ.* **222** 117130.
- BLANGIARDO, M. and CAMELETTI, M. (2015). *Spatial and Spatio-temporal Bayesian Models with R-INLA*. John Wiley & Sons, New York.
- BLEI, D. M., KUCUKELBIR, A. and MCAULIFFE, J. D. (2017). Variational inference: A review for statisticians. *J. Am. Stat. Assoc.* **112** 859–877.
- BOLIN, D., WALLIN, J. and LINDGREN, F. (2019). Latent Gaussian random field mixture models. *Comput. Stat. Data Anal.* **130** 80–93.
- DIEBOLD, F. X. and MARIANO, R. S. (1995). Comparing predictive accuracy. *Journal of Business & economic statistics* **13** 253–263.
- GNEITING, T. (2002). Nonseparable, stationary covariance functions for space-time data. *J. Am. Stat. Assoc.* **97** 590–600.
- HARVEY, D., LEYBOURNE, S. and NEWBOLD, P. (1997). Testing the equality of prediction mean squared errors. *Int. J. Forecast.* **13** 281–291.
- HOUTEKAMER, P. L. and ZHANG, F. (2016). Review of the ensemble Kalman filter for atmospheric data assimilation. *Mon. Weather Rev.* **144** 4489–4532.
- KIRCHGESSNER, P., NERGER, L. and BUNSE-GERSTNER, A. (2014). On the choice of an optimal localization radius in ensemble Kalman filter methods. *Mon. Weather Rev.* **142** 2165–2175.

- LIANG, X., ZOU, T., GUO, B., LI, S., ZHANG, H., ZHANG, S., HUANG, H. and CHEN, S. X. (2015). Assessing Beijing's PM<sub>2.5</sub> pollution: Severity, weather impact, APEC and winter heating. *Proc. R. Soc. A: Math. Phys. Eng. Sci.* **471** 20150257.
- LINDGREN, F., RUE, H. and LINDSTRÖM, J. (2011). An explicit link between Gaussian fields and Gaussian Markov random fields: The stochastic partial differential equation approach. *J. R. Stat. Soc. Ser. B Methodol.* **73** 423–498.
- LINDGREN, F. and RUE, H. (2015). Bayesian spatial modelling with R-INLA. *J. Stat. Softw.* **63** 1–25.
- MITCHELL, H. L., HOUTEKAMER, P. L. and PELLERIN, G. (2002). Ensemble size, balance, and model-error representation in an ensemble Kalman filter. *Mon. Weather Rev.* **130** 2791–2808.
- RODU, J. and KAFADAR, K. (2022). The q-q Boxplot. *J. Comput. Graph. Stat.* **31** 26–39.
- RUE, H. and HELD, L. (2005). *Gaussian Markov Random Fields: Theory and Applications*. Chapman & Hall/CRC, Boca Raton, FL.
- RUE, H., MARTINO, S. and CHOPIN, N. (2009). Approximate Bayesian inference for latent Gaussian models by using integrated nested Laplace approximations. *J. R. Stat. Soc. Ser. B Methodol.* **71** 319–392.
- RUE, H., RIEBLER, A., SØRBYE, S. H., ILLIAN, J. B., SIMPSON, D. P. and LINDGREN, F. K. (2017). Bayesian computing with INLA: a review. *Annu. Rev. Stat. Appl.* **4** 395–421.
- SAHU, S. K., GELFAND, A. E. and HOLLAND, D. M. (2006). Spatio-temporal modeling of fine particulate matter. *J. Agric. Biol. Environ. Stat.* **11** 61–86.
- WOOD, S. N. (2017). *Generalized Additive Models: An Introduction with R*. Chapman & Hall/CRC, New York.
- WOOD, S. N. (2022). Package ‘mgcv’. R package version 1.8-41. <https://CRAN.R-project.org/package=mgcv>.