



# Quantifying the driving factors of particulate matter variabilities in the Beijing-Tianjin-Hebei and Yangtze River Delta regions from 2015 to 2022 by machine learning approach

Zhongfeng Pan<sup>1,2</sup>, Hao Yin<sup>3</sup>, Zhenda Sun<sup>2,4</sup>, Chongyang Li<sup>2,4</sup>, Youwen Sun<sup>1,2</sup>, and Cheng Liu<sup>5,6</sup>

<sup>1</sup>Institutes of Physical Science and Information Technology, Anhui University, Hefei 230601, China

<sup>2</sup>Key Laboratory of Environmental Optics and Technology, Anhui Institute of Optics and Fine Mechanics, HFIPS, Chinese Academy of Sciences, Hefei 230031, China

<sup>3</sup>School of Energy and Environment, City University of Hong Kong, Hong Kong SAR, China

<sup>4</sup>School of Environmental Science and Optoelectronic Technology, University of Science and Technology of China, Hefei 230026, China

<sup>5</sup>Department of Precision Machinery and Precision Instrumentation, University of Science and Technology of China, Hefei 230026, China

<sup>6</sup>Key Laboratory of Precision Scientific Instrumentation of Anhui Higher Education Institutes, University of Science and Technology of China, Hefei 230026, China

**Correspondence:** Hao Yin (haoyin@cityu.edu.hk) and Youwen Sun (ywsun@ahu.edu.cn)

Received: 12 June 2025 – Discussion started: 23 July 2025

Revised: 9 January 2026 – Accepted: 24 January 2026 – Published: 18 February 2026

**Abstract.** Accurately quantifying the relative roles of anthropogenic emissions and meteorological conditions is essential for understanding long term changes in particulate matter (PM). Using ground observations from 40 cities, GEOS-FP meteorology, CEDS emissions, and monthly LightGBM models, this study assesses the drivers of PM<sub>2.5</sub> and PM<sub>10</sub> across the Beijing–Tianjin–Hebei (BTH) and Yangtze River Delta (YRD) regions during 2015–2022. The models demonstrate strong predictive skill ( $R/R^2 = 0.82/0.67$  for PM<sub>2.5</sub> and  $0.81/0.65$  for PM<sub>10</sub>), with consistently high performance across cities. Both pollutants exhibit significant decreasing trends over the study period. Counterfactual experiments show that emission reductions overwhelmingly dominate these improvements. PM<sub>2.5</sub> emission driven changes intensify from  $-9.1 \mu\text{g m}^{-3}$  in 2016 to  $-31.4 \mu\text{g m}^{-3}$  in 2022, while PM<sub>10</sub> reductions strengthen from  $-9.8$  to  $-42.9 \mu\text{g m}^{-3}$ . Meteorology driven contributions appear as positive net anomalies at the interannual scale (approximately  $+2\text{--}4 \mu\text{g m}^{-3}$  for PM<sub>2.5</sub> and  $+0.5\text{--}3 \mu\text{g m}^{-3}$  for PM<sub>10</sub>), indicating that air quality improvements were achieved despite year to year meteorological influences. SHAP attribution highlights 2 m air temperature (T2M), humidity (QV2M), and key precursors as dominant predictors. Interaction diagnostics further indicate that meteorological conditions modulate the effectiveness of precursor emissions, without implying direct causal mechanisms. These results provide a comprehensive data driven assessment of the factors shaping PM evolution in two major urban clusters of China.

## 1 Introduction

Particulate matter (PM) is a significant air pollutant and is also a critical research topic in environmental science due to its diverse sources, complex chemical composition, and profound impacts on human health (Zhang et al., 2022a). Classified by aerodynamic diameter, PM<sub>2.5</sub> (fine particles,  $\leq 2.5 \mu\text{m}$ ) and PM<sub>10</sub> (inhalable particles,  $\leq 10 \mu\text{m}$ ) exert differential impacts on ecosystems and human health owing to their distinct physicochemical properties and environmental behaviors (World Health Organization, 2021). Fine particles (PM<sub>2.5</sub>) penetrate deep into the lungs and cross the alveolar–blood barrier into systemic circulation, while coarser particles (PM<sub>10</sub>) deposit predominantly in the upper respiratory tract (Fu et al., 2024). Chronic exposure to PM<sub>2.5</sub> is linked to respiratory/cardiovascular diseases, declines in lung function, and impairment of the immune system (Franklin et al., 2008; Kioumourtzoglou et al., 2016). Whereas PM<sub>10</sub> aggravates asthma, chronic obstructive pulmonary disease (COPD), and other respiratory conditions (Seaton et al., 1995). As two pivotal economic engines of China, the Beijing-Tianjin-Hebei (BTH) and Yangtze River Delta (YRD) regions, characterized by dense industrial clusters and populations, generate substantial industrial and transportation emissions, with high-intensity production and daily activities resulting in long-standing composite air pollution dominated by PM<sub>2.5</sub>, PM<sub>10</sub>, and ozone (Dai et al., 2021, 2023), posing persistent threats to human health and urban livability. Furthermore, PM pollution acidifies aquatic environments, disrupts ecosystem balance, degrades soils, and contributes to acid rain and terrestrial biosphere damage (Dominici et al., 2014; Jerrett, 2015).

To address severe air pollution problem, Chinese government implemented the Air Pollution Prevention and Control Action Plan (State Council of the People's Republic of China, 2013) and the Three-Year Action Plan for Winning the Blue Sky Defense Battle (State Council of the People's Republic of China, 2018). These initiatives led to substantial reductions in PM concentrations nationwide (Song et al., 2023). However, China's current Ambient Air Quality Standards (GB 3095-2012) stipulate Grade II annual mean limits of  $35 \mu\text{g m}^{-3}$  for PM<sub>2.5</sub> and  $70 \mu\text{g m}^{-3}$  for PM<sub>10</sub>, which significantly exceed the updated WHO guidelines (World Health Organization, 2021).

The dynamics of PM are shaped by anthropogenic precursor emissions – sulfur dioxide (SO<sub>2</sub>), nitrogen oxides (NO<sub>x</sub>), and ammonia (NH<sub>3</sub>) – together with meteorological factors such as temperature, humidity, precipitation, pressure, and wind (Xiao et al., 2021). In addition to these inorganic precursors, volatile organic compounds (VOCs) also play an important role in secondary aerosol formation, particularly through pathways leading to secondary organic aerosols, as recognized in numerous atmospheric chemistry studies. PM<sub>2.5</sub> originates predominantly from traffic and industrial emissions, combustion processes (e.g., cooking,

biomass burning), and secondary formation via atmospheric oxidation to sulfate, nitrate, and organic aerosols (Zhang et al., 2015). PM<sub>10</sub> also includes coarse particles from fugitive dust (construction, agriculture) and secondary coarse-mode particulates (Wu and Huang, 2021). The SO<sub>2</sub>, NO<sub>x</sub>, and NH<sub>3</sub> in the atmosphere can be converted into secondary inorganic aerosols, which significantly regulate PM concentrations (Ding et al., 2019; Feng et al., 2021). Meteorological parameters, such as temperature, relative humidity, precipitation, pressure, and wind, critically influence PM generation, dispersion, and removal (Leung et al., 2018; Zhao et al., 2013). For instance, elevated temperatures accelerate SO<sub>x</sub>/NO<sub>x</sub> oxidation rates and fine PM formation (Chen et al., 2022). High humidity promotes particle hygroscopic growth, gas-to-particle conversion (e.g., secondary organic aerosols), and wet deposition, thereby altering PM size distribution and lifetime. These PM-meteorology interactions exhibit region- and year-specific nonlinear characteristics (Shen et al., 2017), challenging conventional linear modeling approaches (Zhao et al., 2018).

Machine learning (ML), with its capacity to capture complex, nonlinear relationships, has emerged as a powerful tool for atmospheric pollution research (Yin et al., 2022b). ML enhances source apportionment accuracy through multi-source data integration (meteorological, emission, socioeconomic), high-dimensional pattern recognition, and real-time adaptive analysis, enabling identification of complex pollutant interactions (Peng et al., 2024). For PM<sub>2.5</sub> and PM<sub>10</sub> studies, ML facilitates quantitative disentanglement of meteorological and emission contributions, elucidates source-receptor relationships, and informs targeted mitigation strategies.

This study employs the LightGBM framework to quantify the drivers of PM<sub>2.5</sub> and PM<sub>10</sub> variability across the BTH and YRD regions during 2015–2022. By leveraging LightGBM's capability to model nonlinear emission–meteorology–pollution interactions and its efficiency on multi-year, multi-city datasets, the analysis aims to identify the dominant factors governing regional air-quality evolution. The structure of this paper is as follows. Section 2 introduces the datasets used in this study, including national ground-based PM observations (Sect. 2.1), GEOS-FP meteorological reanalysis fields (Sect. 2.2), and the CEDS anthropogenic emission inventory (Sect. 2.3). Section 3 describes the methodological framework. Section 3.1 presents the extraction of city-level meteorological and emission variables from gridded datasets. Section 3.2 outlines the LightGBM modeling workflow, including the feature-engineering strategy (Sect. 3.2.1), the model training and leave-one-year-out cross-validation procedure (Sect. 3.2.2), and the SHAP-based interpretation approach (Sect. 3.2.3). The interannual trend estimation method is detailed in Sect. 3.3, and the counterfactual framework used to separate meteorological and emission contributions is introduced in Sect. 3.4. Section 4 reports the main results, including the interannual

evolution of PM<sub>2.5</sub> and PM<sub>10</sub> (Sect. 4.1), machine-learning model performance and key predictor importance (Sect. 4.2), and the quantified meteorological and emission contributions (Sect. 4.3). Section 5 provides a broader discussion of the identified driving mechanisms in the context of chemical formation pathways and emission-control policies. Finally, Sect. 6 summarizes the key findings and presents implications for future air-quality management and research.

## 2 Data

### 2.1 Observational data from national monitoring sites

The ground-level air pollutant data for the YRD and BTH regions were acquired from the China National Environmental Monitoring Center (CNEMC) network (<https://www.cnemc.cn/>, last access: 31 December 2022), comprising hourly measurements of PM<sub>2.5</sub>, PM<sub>10</sub>, SO<sub>2</sub>, NO<sub>2</sub>, CO, and O<sub>3</sub> concentrations from 2015 to 2022. Observations from multiple monitoring stations within the same city were averaged to derive city-level pollutant concentrations (site-specific details are provided in Table S1 and Fig. S1 in the Supplement). The monitoring network includes 80 stations in the BTH region, covering major cities and areas in Beijing, Tianjin, and Hebei Province, and 197 stations in the YRD region, spanning Shanghai, Jiangsu, Zhejiang, and adjacent provinces. To avoid ambiguity between two cities sharing the same English spelling “Taizhou”, Taizhou City in Jiangsu Province is denoted as TaizhouJS, while Taizhou City in Zhejiang Province is denoted as TaizhouZJ throughout this study.

All national monitoring stations strictly follow the Technical Specifications for Automatic Ambient Air Quality Monitoring (Ministry of Environmental Protection of China, 2013a), ensuring standardized field operation and quality control procedures. City-level PM<sub>2.5</sub> and PM<sub>10</sub> concentrations are released in accordance with the national reference gravimetric method (Ministry of Environmental Protection of China, 2011), which provides the calibration and traceability framework for automated particulate measurements across the monitoring network. Gaseous pollutants (SO<sub>2</sub>, NO<sub>2</sub>, CO, and O<sub>3</sub>) are measured using ultraviolet fluorescence, chemiluminescence, non-dispersive infrared absorption, and ultraviolet photometric analysis, respectively, following the national specifications for continuous gaseous monitoring (Ministry of Environmental Protection of China, 2013b). These unified procedures ensure accuracy, comparability, and long-term stability of pollutant observations.

### 2.2 GEOS-FP meteorological data

Meteorological fields for 2015–2022 were obtained from the GEOS Forward Processing (GEOS-FP) product (<http://geoschemdata.wustl.edu/ExtData/>, last access: 31 December 2020) at a spatial resolution of 0.25° × 0.3125°. GEOS-FP provides hourly assimilated meteorological variables with

relatively high spatial and temporal resolution, enabling refined characterization of mesoscale atmospheric processes over the BTH and YRD regions. Its near-real-time data assimilation framework has been widely applied in regional atmospheric pollution and transport studies, supporting accurate representation of dynamic meteorological conditions (Yin et al., 2021b, 2022a, b). Its near-real-time data assimilation system further improves the accuracy of reanalysis-based meteorological fields and enhances representation of dynamic atmospheric processes (Sun et al., 2021a, b; Yin et al., 2019, 2020, 2021a). The meteorological parameters used in this study include total cloud fraction (CLDTOT), precipitation flux (PRECTOT), 2 m specific humidity (QV2M), 2 m air temperature (T2M), sea-level pressure (SLP), surface downward shortwave radiation (SWGDN), and 10 m zonal (U10M) and meridional (V10M) wind components.

### 2.3 CEDS emission inventory

Anthropogenic emission data for 2015–2022 were obtained from the Community Emissions Data System (CEDS), which provides monthly mean fluxes at a spatial resolution of 0.5° × 0.5°. In this study, we use emissions of CH<sub>2</sub>O, CO, NH<sub>3</sub>, NO<sub>x</sub>, SO<sub>2</sub>, BC, OC, and paraffinic reactive primary emissions (PRPE). Emissions were further categorized into eight sectors: non-combustion agriculture, energy transformation and extraction, industrial combustion and processes, surface transportation, residential and commercial fuel use, solvents, waste treatment and disposal, and international shipping. Anthropogenic emission data for 2015–2022 were derived from the CEDS, a global inventory providing temporally resolved sector-specific emissions. The CEDS framework supports climate change projections and quantifies human-driven interactions between air pollutants and climate systems, critical for assessing health and ecosystem impacts.

## 3 Methodology

### 3.1 Emission and meteorological data extraction

City-level emission and meteorological variables were derived from gridded emission inventories and meteorological reanalysis datasets. For each city, polygon boundaries were obtained from the GADM Level-2 shapefile, and all grid cells whose centers fell within the city polygon were identified.

For emission data, which are expressed as surface fluxes (kg m<sup>-2</sup> s<sup>-1</sup>), the city-scale total emission  $E_{\text{city}}(t)$  was obtained by physical integration over all intersecting grid cells:

$$E_{\text{city}}(t) = \sum_{i,j} F_{ij}(t) A_{ij} r_{ij} \quad (1)$$

where  $F_{ij}(t)$  is the emission flux of grid cell  $(i, j)$  at time  $t$ ;  $A_{ij}$  is the nominal grid-cell area; and  $r_{ij}$  is the fractional overlap between the grid cell and the city polygon. This area-

weighted integration converts surface fluxes into physically consistent city-total emissions.

For meteorological variables, which represent intensive state quantities such as near-surface temperature (T2M), specific humidity (QV2M), and wind components (U10M/V10M), city-level averages were computed as the arithmetic mean of all grid-cell centers located inside the city polygon:

$$X_{\text{city}}(t) = \frac{1}{N} \sum_{i,j} X_{ij}(t) \quad (2)$$

where  $X_{ij}(t)$  is the variable value in cell  $(i, j)$ , and  $N$  is the number of valid grid cells within the city.

This center-based averaging is computationally efficient and provides a representative estimate of the mean meteorological condition. Because the study region lies mainly in the mid-to-low latitudes of China, where grid-cell area variation with latitude is minor, this approximation introduces negligible bias compared with full area-weighted averaging.

Both extraction procedures ensure spatial consistency between emission and meteorological datasets and yield temporally continuous, city-level time series for subsequent modeling.

### 3.2 LightGBM modeling

Light Gradient Boosting Machine (LightGBM) is an efficient and scalable implementation of gradient boosting, extensively applied to regression, classification, and ranking tasks (Yin et al., 2021c). By using a histogram-based decision-tree algorithm, the LightGBM model drastically reduces both computation time and memory usage compared to traditional gradient-boosting methods such as XGBoost (Bian et al., 2023; Zhang et al., 2017). It supports the direct handling of categorical variables without one-hot encoding, improving efficiency when processing high-dimensional datasets. During training, LightGBM grows trees in a leaf-wise (best-first) manner, generating deeper splits along the branch that achieves the largest loss reduction.

In contrast, XGBoost and classical Gradient Boosting Decision Trees (GBDT) use a level-wise growth strategy, which provides stability but can become computationally slower for large-scale data. LightGBM also offers extensive hyperparameter controls, such as maximum tree depth, minimum data in leaf, and feature fraction, to balance model complexity and generalization (Ke et al., 2017). Due to its high predictive accuracy, efficient splitting mechanism, and robust computational capability, LightGBM has become one of the most widely used gradient-boosting frameworks in environmental modeling (Liu et al., 2023; Wang et al., 2022; Zhang et al., 2022b).

Prior to model training, all hourly emission and meteorological inputs were aggregated to monthly means, and the LightGBM models were trained on these monthly-resolved

datasets. This temporal aggregation ensured that all variables shared the same temporal frequency, preventing inconsistencies between hourly and monthly features. It also matched the model input scale with the monthly trend analysis period (2015–2022), thereby avoiding time-scale mismatches between predictor variables and the evaluation framework.

The model performance was evaluated using three widely recognized regression metrics: the correlation coefficient ( $R$ ), indicating the linear relationship between predicted and observed concentrations; the coefficient of determination ( $R^2$ ), measuring the proportion of variance in observations explained by the model; and the root-mean-square error (RMSE), representing the average magnitude of prediction errors. Higher  $R$  and  $R^2$  and lower RMSE indicate stronger predictive ability.

#### 3.2.1 Feature engineering

To ensure interpretability and reduce redundancy, all anthropogenic emissions were aggregated into a single total value for each species (e.g., NOx\_total, SO2\_total). Rather than treating sector-resolved emissions separately, emissions from all sectors were summed into one unified variable per pollutant. This consolidation reduces feature dimensionality, mitigates multicollinearity among sectoral components, and preserves the dominant emission signal, resulting in a more compact and interpretable model structure.

To incorporate temporal information, three seasonal descriptors were introduced: (1) a pair of harmonic terms (month\_sin and month\_cos) representing the cyclic annual progression; (2) a categorical season indicator to reflect broad seasonal regimes. These features provide smooth, physically meaningful representations of seasonality without imposing sharp discontinuities between months.

In addition, two derived emission indicators were constructed to characterize temporal variability and multi-scale coupling:

Seasonal difference (sdiff) represents interannual seasonal changes over a 12 month interval. For each feature  $x$ , the seasonal difference is defined as:

$$\text{sdiff}_t(x) = x_t - x_{t-12} \quad (3)$$

where  $x_t$  denotes the monthly mean value at time  $t$ . This feature highlights annual-scale variations and helps the model capture year-to-year emission variability under identical seasonal conditions.

Rolling detrended residual (detr) represents short-term deviations from a 12 month moving-mean trend. It is expressed as:

$$\text{detr}_t(x) = x_t - \mu_t^{(12)}(x) \quad (4)$$

$$\mu_t^{(12)}(x) = \frac{1}{k_t} \sum_{i=t-11}^t x_i \cdot 1_{\{x_i \text{ valid}\}}, \quad k_t \geq 3 \quad (5)$$

where  $\mu_t^{(12)}(x)$  is the rolling mean computed over the past 12 months, requiring at least three valid data points ( $k_t \geq 3$ ). This feature isolates short-term fluctuations by removing low-frequency seasonal trends.

In practice, the final predictor set comprised 35 variables: eight monthly meteorological parameters, eight species-level aggregated emission totals (e.g., NO<sub>x</sub>\_total, SO<sub>2</sub>\_total, BC\_total), three temporal descriptors (month\_sin, month\_cos, season), and two derived indicators (sdiff and detr) for each emission species. Pollutant concentrations and explicit calendar identifiers (year, month, date) were excluded from the input space to avoid information leakage and to ensure consistent temporal treatment. A complete list of variables used in the model is provided in Table S2.

### 3.2.2 Model training and cross-validation procedure

Before model fitting, all months with missing PM<sub>2.5</sub> or PM<sub>10</sub> observations were removed to ensure consistency between predictors and targets. The derived emission features (sdiff and detr) inherently contain missing values during their initial 12 month computation window; these were addressed using a simple mean-imputation strategy applied within each training fold. Specifically, missing entries for each feature were replaced with the corresponding feature mean computed solely from the training subset of that fold, thereby preventing temporal leakage and ensuring that imputation relied exclusively on information available before prediction.

The LightGBM model was independently trained for each city, pollutant type (PM<sub>2.5</sub> and PM<sub>10</sub>), and cross-validation fold. A leave-one-year-out (LOGO) cross-validation scheme was adopted, whereby data from a single year were held out for testing while the remaining seven years were used for training. This process was repeated sequentially so that each year between 2015 and 2022 served once as the held-out test period. The resulting predictions thus constitute out-of-sample estimates for every year, providing a conservative and temporally robust assessment of model generalizability and avoiding within-year information leakage.

Within each outer leave-one-year-out (LOGO) evaluation split, hyperparameters were tuned using only the training years, with the held-out test year completely excluded from all tuning steps. Specifically, a randomized search was conducted on the training subset to sample 10 candidate configurations from predefined distributions, and candidates were ranked by their cross-validated negative mean squared error (MSE) computed solely within the training data. The selected hyperparameters (learning\_rate, n\_estimators, num\_leaves, max\_depth, min\_child\_samples, subsample, colsample\_bytree, reg\_alpha, and reg\_lambda) were then fixed and used to refit the model on the full training subset of that outer split before generating predictions for

the held-out year. This nested evaluation design ensures that both model fitting and hyperparameter selection rely exclusively on information available prior to prediction, thereby preventing temporal leakage. A consolidated summary of the optimized hyperparameter ranges is presented in Table S3.

### 3.2.3 SHAP interpretation analysis

To interpret the LightGBM outputs and quantify the contribution of individual predictors, the Shapley Additive explanations (SHAP) framework was applied. For each trained model, SHAP values  $s_{i,j}$  represent the marginal contribution of feature  $j$  to the prediction of sample  $i$ . They are derived from cooperative game theory and satisfy the local additivity principle:

$$f(x_i) = E[f(x)] + \sum_{j=1}^m s_{i,j} \quad (6)$$

where  $f(x_i)$  is the model prediction for sample  $i$ , and  $E[f(x)]$  is the expected model output over all samples.

Because this study involves multiple cities across the BTH and YRD regions, a unified regional SHAP ranking was derived by weighting each SHAP importance of city by its corresponding number of valid monthly samples. The weighted importance for feature  $j$  is defined as:

$$I_j^{(w)} = \frac{\sum_{c=1}^N I_{j,c} n_c}{\sum_{c=1}^N n_c} \quad (7)$$

Here  $I_{j,c}$  is the SHAP-based importance of feature  $j$  in city  $c$ , and  $n_c$  is the number of valid samples for that city,  $N$  is the total number of cities.

This weighting strategy ensures that cities with longer and more complete observational records have proportionally greater influence on the regional-level interpretability.

### 3.3 Interannual trend analysis method

To quantify the interannual trends of PM<sub>2.5</sub> and PM<sub>10</sub> concentrations from 2015 to 2022, a linear regression model was employed in this study. For each city, the relationship between annual mean concentration  $y$  and year  $x$  was modeled as:

$$y = \beta_0 + \beta_1 x + \epsilon \quad (8)$$

where  $\beta_0$  represents the intercept (baseline concentration), and  $\epsilon$  denotes the error term. The slope  $\beta_1$ , reflecting the annual rate of concentration change, was estimated via the ordinary least squares (OLS) method. Specifically, the parameters were optimized by minimizing the residual sum of squares (RSS):

$$\arg \min_{\beta_0, \beta_1} \sum_{i=1}^n (y_i - (\beta_0 + \beta_1 x_i))^2 \quad (9)$$

where  $n$  is the sample size (e.g.,  $n = 6$  for the period 2015–2022),  $x_i$  denotes the year, and  $y_i$  represents the corresponding annual mean concentration.

The slope  $\beta_1$  was derived as:

$$\beta_1 = \frac{\text{Cov}(x, y)}{\text{Var}(x)}. \quad (10)$$

The sign of  $\beta_1$  indicates the direction of concentration trends (negative for decreasing, positive for increasing), while its absolute value quantifies the magnitude of change.

### 3.4 Methodology for disentangling meteorological and emission contributions

To quantify meteorological and anthropogenic emission contributions, the trained models were applied in a parallel prediction experiment. For each year from 2016 to 2022, anthropogenic emission features were held fixed at their 2015 levels, while meteorological and temporal predictors were retained at their actual contemporaneous values. This produced a counterfactual concentration series driven solely by meteorological variability. This yielded pollutant concentrations driven solely by meteorological variations (denoted as  $\text{ML}_{2022\text{met}}$  for 2022). The contribution metrics related to 2015 were calculated as follows:

Meteorological contribution ( $\text{ML}_{2022\text{met}}$ ):

$$\text{ML}_{2022\text{met}} = \text{ML}_{15-22} - \text{ML}_{2015}. \quad (11)$$

$\text{ML}_{15-22}$  is the non-emission condition unchanged, the emission condition is fixed as the model prediction result in 2015, and  $\text{ML}_{2015}$  is the model prediction result with unchanged meteorological and emission conditions.

Emission contribution ( $\text{ML}_{2022\text{emis}}$ ):

$$\text{ML}_{2022\text{emis}} = (\text{Obs}_{2022} - \text{Obs}_{2015}) - \text{ML}_{2022\text{met}}. \quad (12)$$

$\text{Obs}_{2022}$  and  $\text{Obs}_{2015}$ : Observed concentrations in 2022 and 2015, respectively.

## 4 Results

### 4.1 Interannual trends of ground-level $\text{PM}_{2.5}$ and $\text{PM}_{10}$

Figure 1 illustrates the interannual evolution of ground-level  $\text{PM}_{2.5}$  and  $\text{PM}_{10}$  concentrations across the BTH and YRD regions from 2015 to 2022, with corresponding annual mean spatial distributions shown in Figs. S2 and S3. Both pollutants exhibit persistent downward trajectories across nearly all cities. The monthly time series and associated OLS trend lines (Fig. S4) further confirm that these declines remain stable throughout the study period, with the most pronounced reductions occurring before 2020. Statistical diagnostics from linear regression (Table S4) indicate that most cities display significant negative trends ( $p < 0.05$ ), underscoring the robustness and consistency of these decreases.

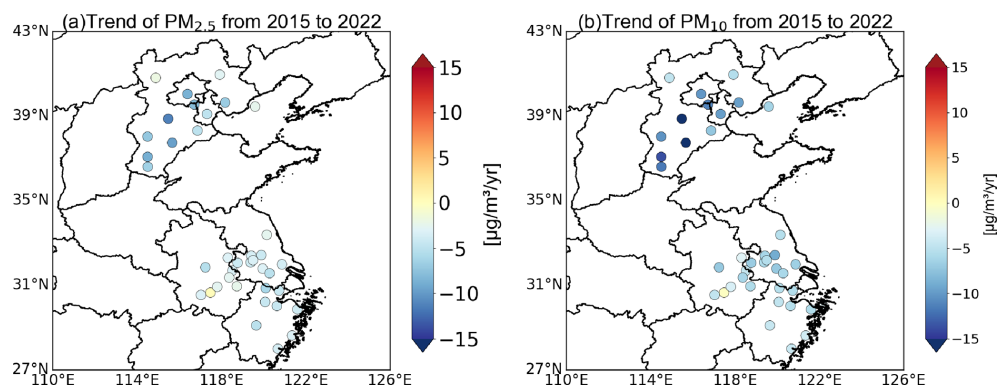
For  $\text{PM}_{2.5}$ , all cities show negative trends, with annual rates ranging from approximately  $-1.53$  to  $-9.46 \mu\text{g m}^{-3} \text{yr}^{-1}$ . Within BTH, the most rapid declines occur in Baoding ( $-9.46 \pm 0.74 \mu\text{g m}^{-3} \text{yr}^{-1}$ ), Hengshui ( $-8.36 \pm 0.67 \mu\text{g m}^{-3} \text{yr}^{-1}$ ), and Xingtai ( $-8.02 \pm 0.56 \mu\text{g m}^{-3} \text{yr}^{-1}$ ), while Chengde ( $-1.98 \pm 0.35 \mu\text{g m}^{-3} \text{yr}^{-1}$ ) and Zhangjiakou ( $-2.00 \pm 0.42 \mu\text{g m}^{-3} \text{yr}^{-1}$ ) exhibit more modest declines. In the YRD region, substantial decreases are observed in Huzhou ( $-5.69 \pm 1.22 \mu\text{g m}^{-3} \text{yr}^{-1}$ ), Hefei ( $-4.83 \pm 0.33 \mu\text{g m}^{-3} \text{yr}^{-1}$ ), and Chuzhou ( $-4.28 \pm 1.12 \mu\text{g m}^{-3} \text{yr}^{-1}$ ), whereas Zhoushan ( $-2.11 \pm 0.14 \mu\text{g m}^{-3} \text{yr}^{-1}$ ) and Taizhou-ZJ ( $-2.72 \pm 0.26 \mu\text{g m}^{-3} \text{yr}^{-1}$ ) experience comparatively smaller reductions. These regional contrasts closely align with the initial concentration levels shown in Fig. S2, where inland BTH cities began with substantially higher baselines (e.g., Baoding and Hengshui exceeded 105 and  $98 \mu\text{g m}^{-3}$  in 2015), enabling larger absolute declines.

For  $\text{PM}_{10}$ , all cities also exhibit significant decreases, with annual trends ranging from roughly  $-1.84$  to  $-13.79 \mu\text{g m}^{-3} \text{yr}^{-1}$ . The largest reductions occur in BTH, particularly Hengshui ( $-13.79 \pm 1.82 \mu\text{g m}^{-3} \text{yr}^{-1}$ ), Baoding ( $-13.07 \pm 1.46 \mu\text{g m}^{-3} \text{yr}^{-1}$ ), and Xingtai ( $-12.76 \pm 1.29 \mu\text{g m}^{-3} \text{yr}^{-1}$ ). In contrast, Chizhou ( $-1.84 \pm 2.03 \mu\text{g m}^{-3} \text{yr}^{-1}$ ) and Zhoushan ( $-2.77 \pm 0.30 \mu\text{g m}^{-3} \text{yr}^{-1}$ ) show the smallest declines, and the  $\text{PM}_{10}$  trend in Chizhou is not statistically significant ( $p = 0.40$ ), consistent with its minimal decline rate. As with  $\text{PM}_{2.5}$ , cities starting with the highest  $\text{PM}_{10}$  levels – such as Baoding ( $174.6 \mu\text{g m}^{-3}$  in 2015) and Hengshui ( $175.9 \mu\text{g m}^{-3}$ ) – exhibit the steepest reductions, whereas cities with lower initial levels (e.g., Zhoushan at  $\sim 46.8 \mu\text{g m}^{-3}$ ) show correspondingly smaller absolute decreases.

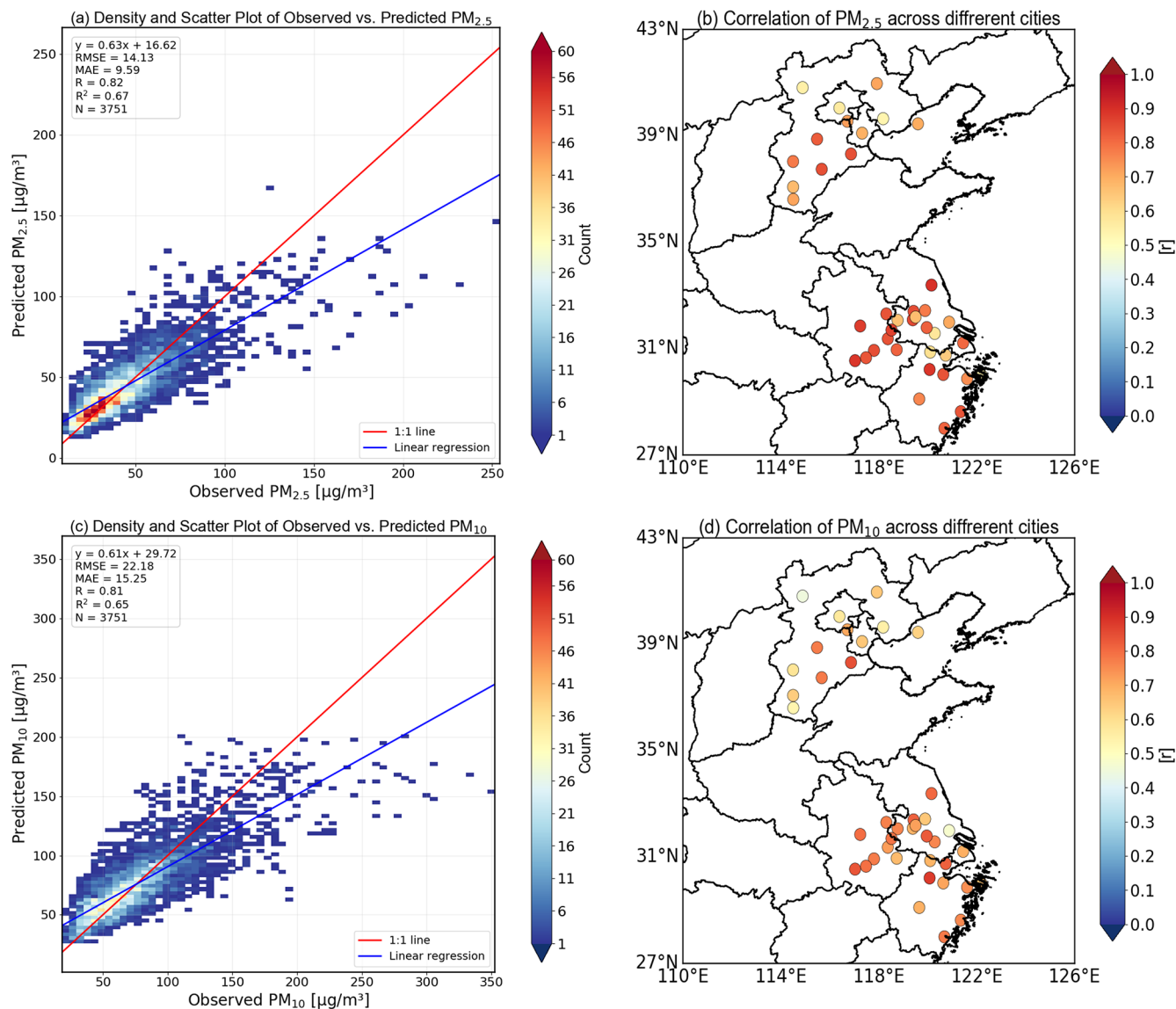
Overall, both  $\text{PM}_{2.5}$  and  $\text{PM}_{10}$  decrease more rapidly in BTH than in YRD, reflecting regional differences in initial emissions, industrial structure, and the intensity of mitigation policies. The widespread statistical significance of the trends (Table S4) supports the conclusion that both regions experienced sustained and robust improvements in air quality during 2015–2022.

### 4.2 Machine learning model performance and variable importance

Figure 2 summarizes the predictive performance of the LightGBM models for  $\text{PM}_{2.5}$  and  $\text{PM}_{10}$  across all cities. The density scatterplots (Fig. 2a and c) show good agreement between predictions and observations, with overall  $R/R^2$  values of 0.82/0.67 for  $\text{PM}_{2.5}$  and 0.81/0.65 for  $\text{PM}_{10}$ , respectively. City-level performance (Fig. 2b and d) further indicates that the majority of cities achieve correlation coefficients exceeding 0.70, as detailed in Table S5. While the predictive skill varies among cities, particularly for  $\text{PM}_{10}$ , such variability is expected given the stronger influence of

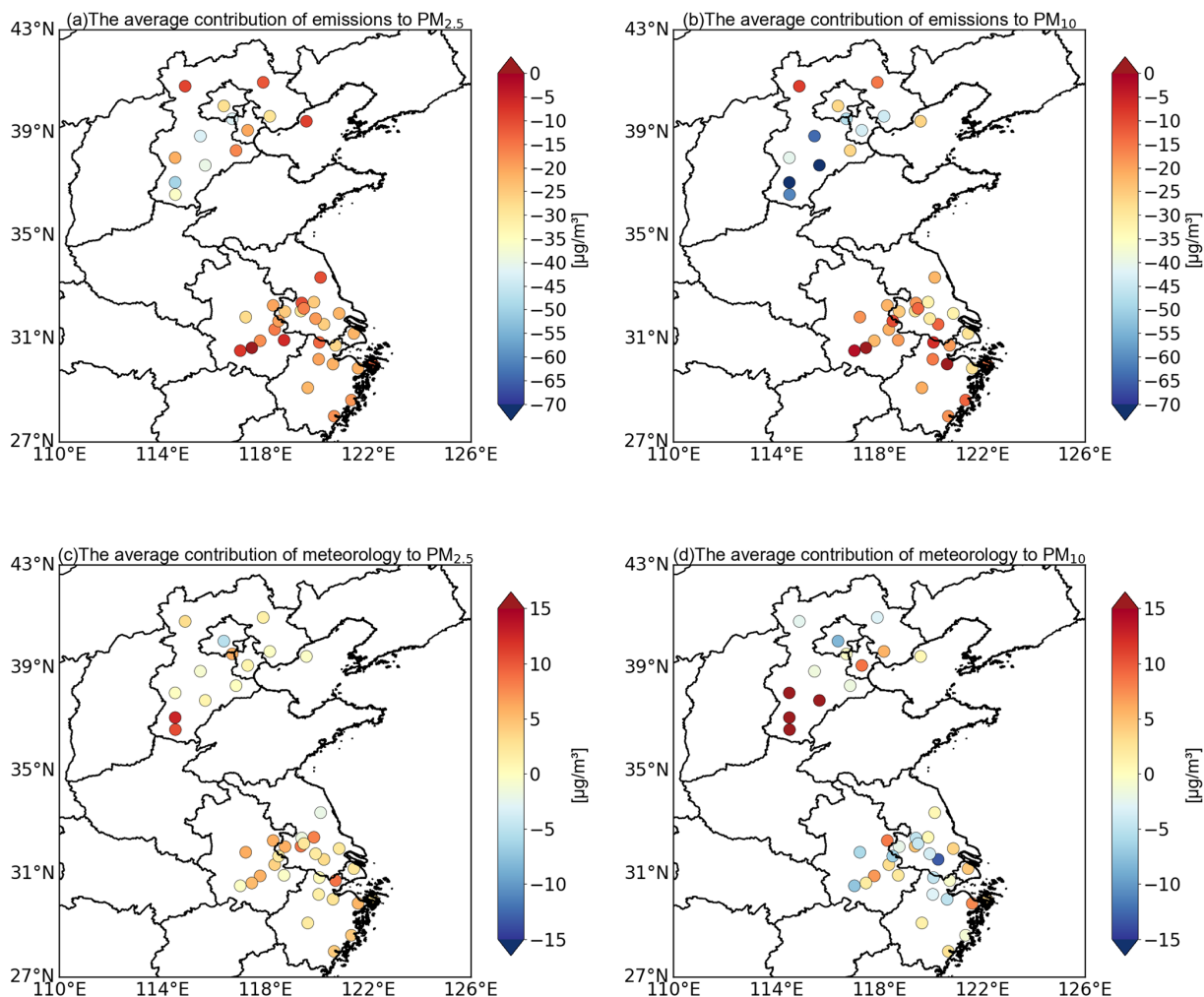


**Figure 1.** Interannual variation trends of  $\text{PM}_{2.5}$  (a) and  $\text{PM}_{10}$  (b) in each city during 2015–2022.



**Figure 2.** The density scatter plots of  $\text{PM}_{2.5}$  (a) and  $\text{PM}_{10}$  (c) concentrations observed and predicted, respectively. The correlation of  $\text{PM}_{2.5}$  (b) and  $\text{PM}_{10}$  (d) in each city over BTH and YRD regions, respectively.



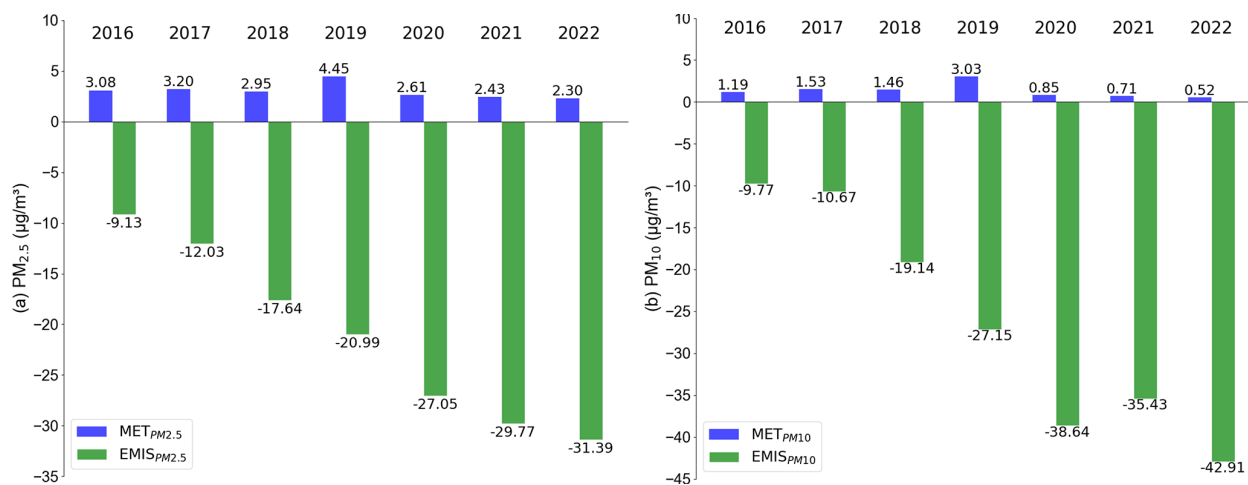


**Figure 4.** The average contributions of emissions and meteorological variables to  $\text{PM}_{2.5}$  (for **a** and **c**) and  $\text{PM}_{10}$  (for **b** and **d**), respectively.

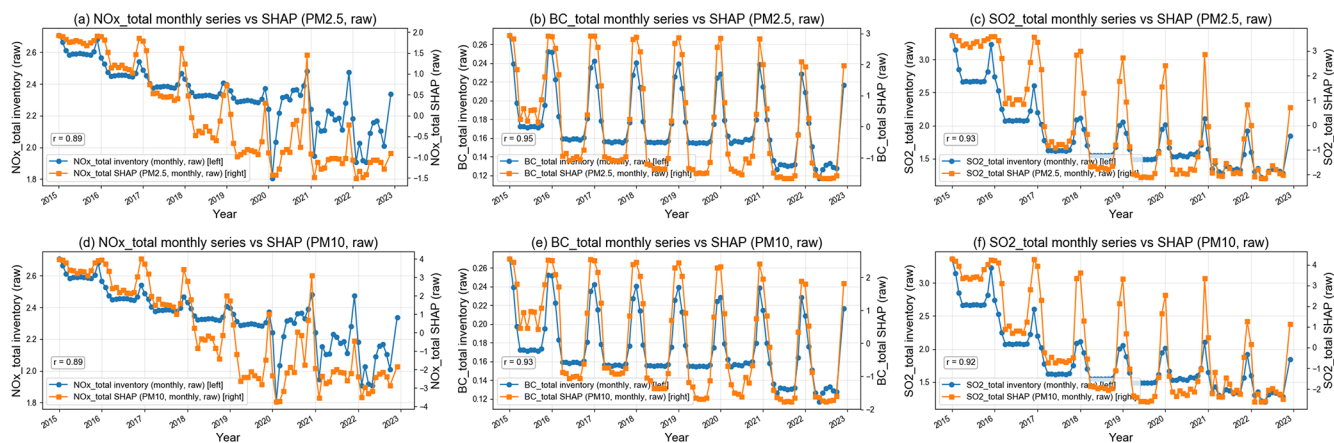
from  $-9.13 \mu\text{g m}^{-3}$  in 2016 to  $-31.39 \mu\text{g m}^{-3}$  in 2022, while  $\text{PM}_{10}$  reductions strengthen from  $-9.77$  to  $-42.91 \mu\text{g m}^{-3}$  over the same period. These progressively stronger negative contributions are consistent with sustained decreases in major emission species, including  $\text{NO}_x_{\text{total}}$ ,  $\text{SO}_2_{\text{total}}$ , and  $\text{BC}_{\text{total}}$ . In contrast, meteorological contributions remain positive throughout all years, on the order of approximately  $2\text{--}4 \mu\text{g m}^{-3}$  for  $\text{PM}_{2.5}$  and  $0.5\text{--}3 \mu\text{g m}^{-3}$  for  $\text{PM}_{10}$ , indicating that interannual meteorological conditions, on average, tended to offset part of the emission-driven reductions rather than reinforce them. Although  $\text{PM}_{10}$  is often considered more responsive to short-term meteorological variability, the smaller meteorology-driven contributions for  $\text{PM}_{10}$  in Fig. 5 should be interpreted in the context of interannual net effects. Meteorological influences on  $\text{PM}_{10}$  often involve multiple processes that can act in opposite directions, such as enhanced resuspension versus enhanced dispersion and removal, leading to partial cancellation when aggregated at annual timescales. By contrast, meteorological effects on

$\text{PM}_{2.5}$ , particularly those related to temperature and humidity, tend to influence secondary formation and particle growth in a more directionally consistent manner at the interannual scale. This consistency allows their impacts to accumulate into a clearer net contribution. As a result, the relative magnitude of meteorology-driven changes appears larger for  $\text{PM}_{2.5}$  than for  $\text{PM}_{10}$  in Fig. 5, without implying a weaker overall meteorological sensitivity of  $\text{PM}_{10}$ .

Temporal relationships between individual emission indicators and their contributions are shown in Fig. 6. For  $\text{PM}_{2.5}$ , month-to-month changes in  $\text{NO}_x_{\text{total}}$ ,  $\text{SO}_2_{\text{total}}$ , and  $\text{BC}_{\text{total}}$  are strongly correlated with their corresponding emission-driven contributions ( $R = 0.89\text{--}0.95$ ). A similar correspondence is observed for  $\text{PM}_{10}$  ( $R = 0.89\text{--}0.93$ ). These strong associations indicate that the estimated emission contributions respond consistently to changes in emission levels, with higher emissions generally linked to larger positive contributions and lower emissions associated with weaker contributions. This relationship persists despite the



**Figure 5.** The averaging the emission or meteorological contributions to PM<sub>2.5</sub> (a) and PM<sub>10</sub> (b) of each year relative to 2015.



**Figure 6.** Monthly emission inventories of NO<sub>x</sub><sub>total</sub>, BC<sub>total</sub>, and SO<sub>2</sub><sub>total</sub> (left y axis) and their corresponding SHAP values (right y axis) for PM<sub>2.5</sub> and PM<sub>10</sub> from 2015 to 2022. Panels (a)–(c) show NO<sub>x</sub><sub>total</sub>, BC<sub>total</sub>, and SO<sub>2</sub><sub>total</sub> and their SHAP contributions to PM<sub>2.5</sub>, while panels (d)–(f) present the corresponding relationships for PM<sub>10</sub>. Emission inventories and SHAP values are plotted on separate y axes to account for their different magnitudes.

nonlinear structure of the model, which adjusts contribution magnitudes according to the prevailing meteorological conditions.

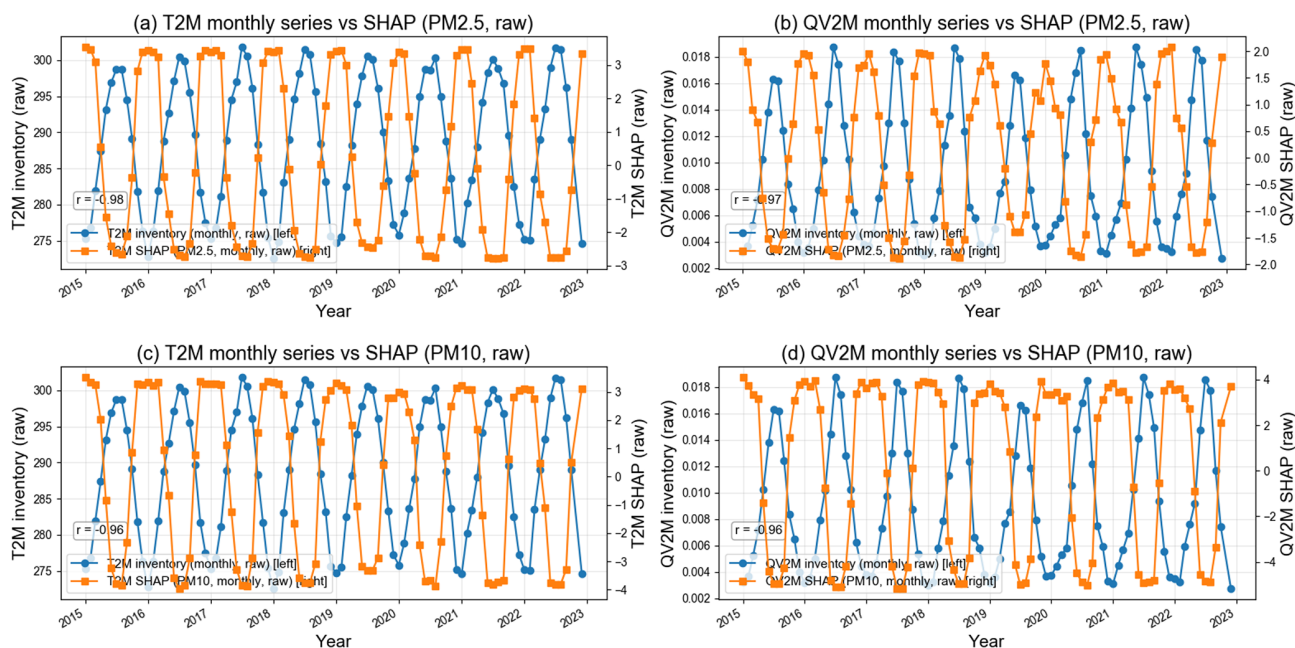
In contrast, meteorological drivers show a different temporal behavior (Fig. 7). Temperature and humidity exhibit strong negative correlations with their meteorology-driven contributions for both pollutants ( $R \approx -0.96$  to  $-0.98$ ). Importantly, these correlations describe relative contributions at the interannual scale after emission-driven trends have been removed, rather than direct or instantaneous responses of particulate matter to meteorological forcing. At this aggregated temporal scale, years with higher temperature or moisture do not necessarily correspond to larger net meteorological contributions, because multiple meteorological influences with opposing effects can partially offset each other. These results therefore point to a scale-dependent influence of meteorology, rather than contradicting process-based studies that em-

phasize the role of temperature and humidity in promoting particle formation at shorter timescales.

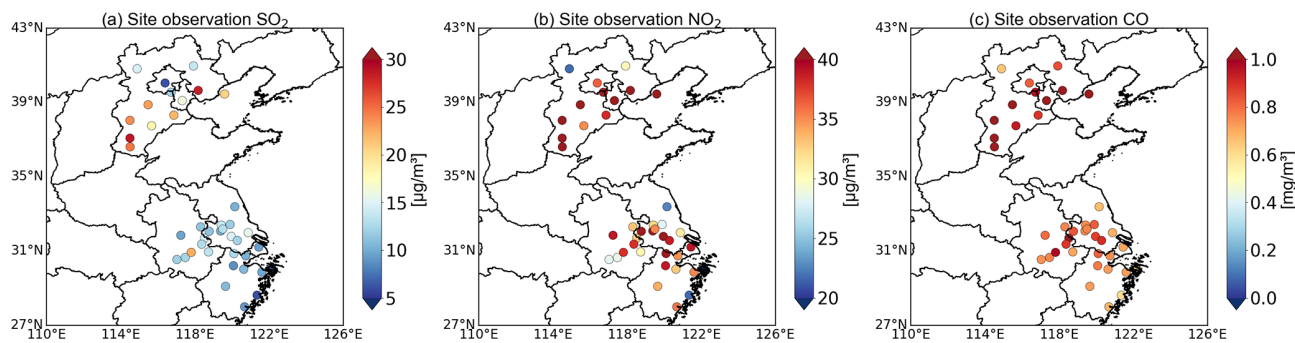
Taken together, the increasingly negative emission-driven contributions shown in Fig. 5, combined with relatively modest and predominantly positive net meteorology-driven anomalies at the interannual scale, indicate that the observed improvements in PM<sub>2.5</sub> and PM<sub>10</sub> during 2016–2022 were mainly driven by sustained emission reductions. Meteorological conditions influenced year-to-year fluctuations but did not reverse the overall downward trends, highlighting the dominant role of emission control measures in shaping the long-term evolution of both pollutants.

## 5 Discussions

To provide additional context for the SHAP-derived attributions, we examine the spatial distributions and statistical



**Figure 7.** Monthly meteorological variables (T2M and QV2M; left y axis) and their corresponding SHAP values (right y axis) for  $\text{PM}_{2.5}$  and  $\text{PM}_{10}$  from 2015 to 2022. Panels (a) and (b) show the relationships for  $\text{PM}_{2.5}$ , while panels (c) and (d) present the corresponding results for  $\text{PM}_{10}$ . Separate y axes are used to account for differences in magnitude between meteorological variables and SHAP values.



**Figure 8.** The average concentrations of  $\text{SO}_2$  (a),  $\text{NO}_2$  (b),  $\text{CO}$  (c), respectively, during 2015 to 2020 over BTH and YRD regions.

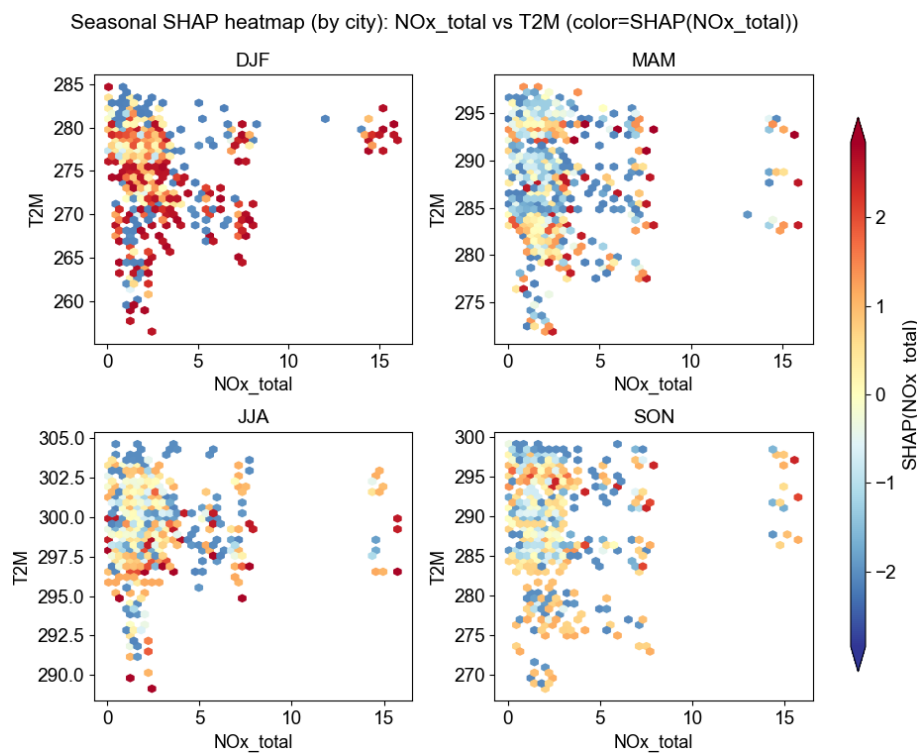
associations of key precursor species. Figure 8 shows that ambient concentrations of  $\text{CO}$ ,  $\text{NO}_2$ , and  $\text{SO}_2$  display pronounced spatial contrasts across the BTH and YRD regions. Heavy industrial cities such as Tangshan, Xingtai, Handan, and Baoding consistently exhibit higher levels of these pollutants, whereas several coastal YRD cities (e.g., Zhoushan and TaizhouZJ) show substantially lower concentrations, reflecting differences in energy structure and industrial activity. These spatial patterns offer a useful background for interpreting the correlation results summarized in Table 1.

Across all cities,  $\text{CO}$  shows the strongest statistical association with  $\text{PM}_{2.5}$  (mean  $R \approx 0.72$ ) and  $\text{PM}_{10}$  (mean  $R \approx 0.62$ ), followed by  $\text{NO}_2$  ( $R \approx 0.59/0.57$ ) and  $\text{SO}_2$  ( $R \approx 0.49/0.48$ ). The relatively strong  $\text{CO}$ – $\text{PM}$  relationships are consistent with their shared combustion-related origins, including traffic emissions, industrial fuel use, and residential

**Table 1.** The correlation values among  $\text{SO}_2$ ,  $\text{NO}_2$ ,  $\text{CO}$  and  $\text{PM}_{2.5}/\text{PM}_{10}$ , respectively.

	$\text{SO}_2$	$\text{NO}_2$	$\text{CO}$
$\text{PM}_{2.5}$	0.49	0.59	0.72
$\text{PM}_{10}$	0.48	0.57	0.62

heating, and align with recent multi-platform and top-down studies reporting tight coupling among  $\text{CO}$ ,  $\text{NO}_x$ , carbonaceous aerosols, and combustion-related  $\text{PM}_{2.5}$  (Tiwari et al., 2025; Wang et al., 2021, 2025). In contrast, the weaker correlations involving  $\text{SO}_2$ , despite elevated concentrations in several northern cities, likely reflect the long-term effectiveness of desulfurization policies, which have reduced sulfate



**Figure 9.** Seasonal scatter plots of NO<sub>x</sub>\_total versus T2M across all cities, with point colors representing the SHAP contribution of NO<sub>x</sub>\_total. The four panels correspond to the DJF, MAM, JJA, and SON seasons.

formation and altered secondary aerosol composition. NO<sub>2</sub> correlations fall between those of CO and SO<sub>2</sub>, indicating sustained contributions from traffic and industrial sources while also reflecting evolving chemical pathways and emission controls. Building on this context, Fig. 9 illustrates the joint dependence of NO<sub>x</sub>\_total-related SHAP contributions on temperature for PM<sub>2.5</sub>. Lower temperatures are associated with stronger NO<sub>x</sub>-related contributions, whereas this influence weakens at higher temperatures. This pattern is consistent with established seasonal behavior in which cold conditions favor the persistence of nitrate-related particulate matter and reduced atmospheric mixing, while warmer conditions limit nitrate effectiveness and enhance dilution. Here, the interaction plot is intended as a diagnostic illustration showing how SHAP-derived contributions vary across the observed temperature range, rather than as a standalone mechanistic attribution. We do not present analogous interaction panels for PM<sub>10</sub> or for humidity-related variables. Coarse particles are more strongly influenced by mechanically driven processes such as dust resuspension and surface conditions, leading to less systematic chemical responses to temperature. Similarly, humidity affects multiple competing processes, including hygroscopic growth, aqueous-phase reactions, boundary-layer suppression, and wet removal, making it difficult to isolate a single, interpretable interaction at the seasonal scale.

Taken together, the analyses in this section are intended to complement the main SHAP-based attribution results by providing empirical context and internal consistency checks. They suggest that the statistical relationships identified by the model are broadly compatible with known emission structures, chemical regimes, and seasonal behavior, while reinforcing that the primary physical interpretation of PM variability is established by the feature importance rankings and contribution analyses discussed in earlier sections.

## 6 Conclusions

This study integrates multi-source emission inventories, GEOS-FP meteorological fields, and ground-based observations to investigate the drivers of PM<sub>2.5</sub> and PM<sub>10</sub> variability across the BTH and YRD regions during 2015–2022 using a unified LightGBM modeling and attribution framework. Both pollutants exhibit clear and statistically significant declines across most cities, supported by strong model performance under a rigorous leave-one-year-out cross-validation design ( $R/R^2 = 0.82/0.67$  for PM<sub>2.5</sub>;  $0.81/0.65$  for PM<sub>10</sub>). The consistent agreement between predictions and observations across regions and years demonstrates the robustness of the machine-learning-based representation of emission–meteorology–PM relationships at the interannual scale.

A central contribution of this work is the integration of SHAP-based interpretability diagnostics, which enables

transparent attribution while retaining the flexibility of non-linear learning. The attribution analyses indicate that anthropogenic emission reductions are the primary driver of the observed PM improvements. For  $\text{PM}_{2.5}$ , emission-driven decreases intensify from  $-9.13 \mu\text{g m}^{-3}$  in 2016 to  $-31.39 \mu\text{g m}^{-3}$  in 2022, while corresponding reductions for  $\text{PM}_{10}$  strengthen from  $-9.77$  to  $-42.91 \mu\text{g m}^{-3}$ . These progressively larger negative contributions are consistent with sustained declines in major precursors, including  $\text{NO}_x$ ,  $\text{SO}_2$ , and BC. In contrast, meteorology-driven contributions appear as comparatively smaller net anomalies at the interannual scale, indicating that year-to-year meteorological variability modulated PM levels but did not offset the long-term emission-driven downward trends.

The consistency between precursor emissions and their SHAP-derived contributions at the monthly scale, together with scale-aware temperature and humidity patterns, provides additional confidence that the attribution results are compatible with established emission structures, seasonal behavior, and chemical regimes. Rather than serving as direct mechanistic proofs, these diagnostics offer internal consistency checks that support the interpretability of the machine-learning framework.

Several methodological strengths emerge from this study. The use of harmonized monthly emission totals, multi-scale temporal descriptors (sdiff and detr), and SHAP-based interaction diagnostics facilitates a physically informed interpretation of complex emission–meteorology coupling. In addition, the explicit construction of counterfactual predictions provides a transparent approach for separating meteorological and anthropogenic influences, offering a reproducible pathway for applying machine-learning-based attribution methods to other regions and pollutants.

Despite these advantages, some limitations remain. The monthly temporal resolution cannot fully capture short-term meteorological or chemical processes; uncertainties in bottom-up emission inventories and meteorological reanalyses may affect the absolute magnitudes of estimated contributions; and, as with all data-driven approaches, causal relationships cannot be inferred directly from statistical associations. Future work will extend this framework to explicitly examine coarse-mode particles by analyzing  $\text{PM}_{2.5-10}$ , enabling a clearer separation of fine particle and coarse particle drivers. Additional developments may include higher-frequency observations, expanded precursor coverage (e.g., VOCs), and hybrid machine learning and chemical transport modeling to further improve process interpretability and physical fidelity.

Overall, this study demonstrates that sustained multi-sector anthropogenic emission reductions, rather than meteorological variability, primarily explain the observed decreases in  $\text{PM}_{2.5}$  and  $\text{PM}_{10}$  from 2015 to 2022. The results highlight the value of physically informed machine learning tools, coupled with SHAP-based interpretability, for diagnos-

ing long term air quality evolution and supporting emission control strategies grounded in observational evidence.

**Code and data availability.** The code and data for this study can be found on <https://doi.org/10.5281/zenodo.17779780> (Pan et al., 2025).

**Supplement.** The supplement related to this article is available online at <https://doi.org/10.5194/acp-26-2545-2026-supplement>.

**Author contributions.** HY and YS designed this study. ZP wrote the paper with help from HY and YS. ZP contributed to analysis of the data for this study. All co-authors commented on this study.

**Competing interests.** The contact author has declared that none of the authors has any competing interests.

**Disclaimer.** Publisher's note: Copernicus Publications remains neutral with regard to jurisdictional claims made in the text, published maps, institutional affiliations, or any other geographical representation in this paper. The authors bear the ultimate responsibility for providing appropriate place names. Views expressed in the text are those of the authors and do not necessarily reflect the views of the publisher.

**Acknowledgements.** This work is jointly supported by Excellent Young Scientists Fund of the National Natural Science Foundation of China (62322514), Anhui Science Fund for Distinguished Young Scholars (2308085J25), and National Key Research and Development Program of China (2023YFC3709502, 2022YFC3700100).

**Financial support.** This research has been supported by the Excellent Young Scientists Fund of the National Natural Science Foundation of China (62322514), the Science Fund for Distinguished Young Scholars of Anhui Province (2308085J25), and National Key Research and Development Program of China (2023YFC3709502, 2022YFC3700100).

**Review statement.** This paper was edited by Jason Cohen and reviewed by four anonymous referees.

## References

- Bian, L., Qin, X., Zhang, C., Guo, P., and Wu, H.: Application, interpretability and prediction of machine learning method combined with LSTM and LightGBM—a case study for runoff simulation in an arid area, *J. Hydrol.*, 625, 130091, <https://doi.org/10.1016/j.jhydrol.2023.130091>, 2023.

- Chen, Y., Su, W., Xing, C., Yin, H., Lin, H., Zhang, C., Liu, H., Hu, Q., and Liu, C.: Kilometer-level glyoxal retrieval via satellite for anthropogenic volatile organic compound emission source and secondary organic aerosol formation identification, *Remote Sens. Environ.*, 270, 112852, <https://doi.org/10.1016/j.rse.2021.112852>, 2022.
- Dai, H., Zhu, J., Liao, H., Li, J., Liang, M., Yang, Y., and Yue, X.: Co-occurrence of ozone and PM<sub>2.5</sub> pollution in the Yangtze River Delta over 2013–2019: Spatiotemporal distribution and meteorological conditions, *Atmospheric Res.*, 249, 105363, <https://doi.org/10.1016/j.atmosres.2020.105363>, 2021.
- Dai, H., Liao, H., Li, K., Yue, X., Yang, Y., Zhu, J., Jin, J., Li, B., and Jiang, X.: Compositing analyses of the chemical and physical characteristics of co-polluted days by ozone and PM<sub>2.5</sub> over 2013–2020 in the Beijing–Tianjin–Hebei region, *Atmos. Chem. Phys.*, 23, 23–39, <https://doi.org/10.5194/acp-23-23-2023>, 2023.
- Ding, A., Huang, X., Nie, W., Chi, X., Xu, Z., Zheng, L., Xu, Z., Xie, Y., Qi, X., Shen, Y., Sun, P., Wang, J., Wang, L., Sun, J., Yang, X.-Q., Qin, W., Zhang, X., Cheng, W., Liu, W., Pan, L., and Fu, C.: Significant reduction of PM<sub>2.5</sub> in eastern China due to regional-scale emission control: evidence from SORPES in 2011–2018, *Atmos. Chem. Phys.*, 19, 11791–11801, <https://doi.org/10.5194/acp-19-11791-2019>, 2019.
- Dominici, F., Greenstone, M., and Sunstein, C. R.: Particulate Matter Matters, *Science*, 344, 257–259, <https://doi.org/10.1126/science.1247348>, 2014.
- Feng, X., Tian, Y., Xue, Q., Song, D., Huang, F., and Feng, Y.: Measurement report: Spatiotemporal and policy-related variations of PM<sub>2.5</sub> composition and sources during 2015–2019 at multiple sites in a Chinese megacity, *Atmos. Chem. Phys.*, 21, 16219–16235, <https://doi.org/10.5194/acp-21-16219-2021>, 2021.
- Franklin, M., Koutrakis, P., and Schwartz, P.: The role of particle composition on the association between PM<sub>2.5</sub> and mortality, *Epidemiology*, 19, 680–689, <https://doi.org/10.1097/ede.0b013e3181812bb7>, 2008.
- Fu, L., Guo, Y., Zhu, Q., Chen, Z., Yu, S., Xu, J., Tang, W., Wu, C., He, G., Hu, J., Zeng, F., Dong, X., Yang, P., Lin, Z., Wu, F., Liu, T., and Ma, W.: Effects of long-term exposure to ambient fine particulate matter and its specific components on blood pressure and hypertension incidence, *Environ. Int.*, 184, 108464, <https://doi.org/10.1016/j.envint.2024.108464>, 2024.
- Jerrett, M.: Atmospheric science: The death toll from air-pollution sources, *Nature*, 525, 330–331, <https://doi.org/10.1038/525330a>, 2015.
- Ke, G., Meng, Q., Finley, T., Wang, T., Chen, W., Ma, W., Ye, Q., and Liu, T.-Y.: LightGBM: A Highly Efficient Gradient Boosting Decision Tree, *Advances in Neural Information Processing Systems*, in: NIPS'17: Proceedings of the 31st International Conference on Neural Information Processing Systems, <https://dl.acm.org/doi/10.5555/3294996.3295074>, 2017.
- Kioumourtzoglou, M.-A., Schwartz, J. D., Weisskopf, M. G., Melly, S. J., Wang, Y., Dominici, F., and Zanobetti, A.: Long-term PM<sub>2.5</sub> Exposure and Neurological Hospital Admissions in the Northeastern United States, *Environ. Health Perspect.*, 124, 23–29, <https://doi.org/10.1289/ehp.1408973>, 2016.
- Leung, D. M., Tai, A. P. K., Mickley, L. J., Moch, J. M., van Donkelaar, A., Shen, L., and Martin, R. V.: Synoptic meteorological modes of variability for fine particulate matter (PM<sub>2.5</sub>) air quality in major metropolitan regions of China, *Atmos. Chem. Phys.*, 18, 6733–6748, <https://doi.org/10.5194/acp-18-6733-2018>, 2018.
- Liu, C., Yang, Y., Wang, H., Ren, L., Wei, J., Wang, P., and Liao, H.: Influence of Spatial Dipole Pattern in Asian Aerosol Changes on East Asian Summer Monsoon, *J. Clim.*, 36, 1575–1585, <https://doi.org/10.1175/JCLI-D-22-0335.1>, 2023.
- Ministry of Environmental Protection of China: Determination of Atmospheric Particulates PM<sub>10</sub> and PM<sub>2.5</sub> in Ambient Air by Gravimetric Method, Ministry of Environmental Protection of China, [https://english.mee.gov.cn/Resources/standards/Air\\_Environment/air\\_method/201111/t20111101\\_219390.shtml](https://english.mee.gov.cn/Resources/standards/Air_Environment/air_method/201111/t20111101_219390.shtml) (last access: 8 February 2026), 2011.
- Ministry of Environmental Protection of China: Specifications and Test Procedures for Ambient Air Quality Continuous Automated Monitoring System for PM<sub>10</sub> and PM<sub>2.5</sub>, Ministry of Environmental Protection of China, [https://english.mee.gov.cn/Resources/standards/Air\\_Environment/air\\_method/201308/t20130816\\_257558.shtml](https://english.mee.gov.cn/Resources/standards/Air_Environment/air_method/201308/t20130816_257558.shtml) (last access: 8 February 2026), 2013a.
- Ministry of Environmental Protection of China: Specifications and Test Procedures for Ambient Air Quality Continuous Automated Monitoring System for SO<sub>2</sub>, NO<sub>2</sub>, O<sub>3</sub> and CO, Ministry of Environmental Protection of China, [https://english.mee.gov.cn/Resources/standards/Air\\_Environment/air\\_method/201308/t20130816\\_257557.shtml](https://english.mee.gov.cn/Resources/standards/Air_Environment/air_method/201308/t20130816_257557.shtml) (last access: 8 February 2026), 2013b.
- Pan, Z., Sun, Y., Yin, H., and Liu, C.: Quantifying the driving factors of particulate matter variabilities in the Beijing-Tianjin-Hebei and Yangtze River Delta regions from 2015 to 2022 by machine learning approach, Version v2, Zenodo [data set/code], <https://doi.org/10.5281/zenodo.17779780>, 2025.
- Peng, Z., Zhang, B., Wang, D., Niu, X., Sun, J., Xu, H., Cao, J., and Shen, Z.: Application of machine learning in atmospheric pollution research: A state-of-art review, *Sci. Total Environ.*, 910, 168588, <https://doi.org/10.1016/j.scitotenv.2023.168588>, 2024.
- Seaton, A., MacNee, W., Donaldson, K., and Godden, D.: Particulate air pollution and acute health effects, *Lancet Lond. Engl.*, 345, 176–178, [https://doi.org/10.1016/s0140-6736\(95\)90173-6](https://doi.org/10.1016/s0140-6736(95)90173-6), 1995.
- Shen, L., Mickley, L. J., and Murray, L. T.: Influence of 2000–2050 climate change on particulate matter in the United States: results from a new statistical model, *Atmos. Chem. Phys.*, 17, 4355–4367, <https://doi.org/10.5194/acp-17-4355-2017>, 2017.
- Song, X.-H., Yan, L., Liu, W., He, J.-Y., Wang, Y.-C., Huang, T.-L., Li, Y.-Y., Chen, M., Meng, J.-J., and Hou, Z.-F.: Spatiotemporal Distribution Characteristics of Co-pollution of PM<sub>2.5</sub> and Ozone over BTH with Surrounding Area from 2015 to 2021, *Huan Jing Ke Xue Huanjing Kexue*, 44, 1841–1851, <https://doi.org/10.13227/j.hjhx.202205089>, 2023.
- State Council of the People's Republic of China: Air Pollution Prevention and Control Action Plan, State Council of the People's Republic of China, [https://english.mee.gov.cn/News\\_service/infocus/201309/t20130924\\_260707.shtml](https://english.mee.gov.cn/News_service/infocus/201309/t20130924_260707.shtml) (last access: 8 February 2026), 2013.
- State Council of the People's Republic of China: Three-Year Action Plan for Winning the Blue Sky Defense Battle, State Council of the People's Republic of China, [https://english.mee.gov.cn/News\\_service/news\\_release/201807/t20180713\\_446624.shtml](https://english.mee.gov.cn/News_service/news_release/201807/t20180713_446624.shtml) (last access: 8 February 2026), 2018.

- Sun, Y., Yin, H., Liu, C., Zhang, L., Cheng, Y., Palm, M., Notholt, J., Lu, X., Vigouroux, C., Zheng, B., Wang, W., Jones, N., Shan, C., Qin, M., Tian, Y., Hu, Q., Meng, F., and Liu, J.: Mapping the drivers of formaldehyde (HCHO) variability from 2015 to 2019 over eastern China: insights from Fourier transform infrared observation and GEOS-Chem model simulation, *Atmos. Chem. Phys.*, 21, 6365–6387, <https://doi.org/10.5194/acp-21-6365-2021>, 2021a.
- Sun, Y., Yin, H., Liu, C., Mahieu, E., Notholt, J., Té, Y., Lu, X., Palm, M., Wang, W., Shan, C., Hu, Q., Qin, M., Tian, Y., and Zheng, B.: The reduction in C<sub>2</sub>H<sub>6</sub> from 2015 to 2020 over Hefei, eastern China, points to air quality improvement in China, *Atmos. Chem. Phys.*, 21, 11759–11779, <https://doi.org/10.5194/acp-21-11759-2021>, 2021b.
- Tiwari, P., Cohen, J. B., Lu, L., Wang, S., Li, X., Guan, L., Liu, Z., Li, Z., and Qin, K.: Multi-platform observations and constraints reveal overlooked urban sources of black carbon in Xuzhou and Dhaka, *Commun. Earth Environ.*, 6, 38, <https://doi.org/10.1038/s43247-025-02012-x>, 2025.
- Wang, S., Cohen, J. B., Deng, W., Qin, K., and Guo, J.: Using a New Top-Down Constrained Emissions Inventory to Attribute the Previously Unknown Source of Extreme Aerosol Loadings Observed Annually in the Monsoon Asia Free Troposphere, *Earths Future*, 9, e2021EF002167, <https://doi.org/10.1029/2021EF002167>, 2021.
- Wang, S., Cohen, J. B., Guan, L., Lu, L., Tiwari, P., and Qin, K.: Observationally constrained global NO<sub>x</sub> and CO emissions variability reveals sources which contribute significantly to CO<sub>2</sub> emissions, *Npj Clim. Atmos. Sci.*, 8, 87, <https://doi.org/10.1038/s41612-025-00977-2>, 2025.
- Wang, X., Xue, Y., Jin, C., Sun, Y., and Li, N.: Spatial downscaling of surface ozone concentration calculation from remotely sensed data based on mutual information, *Front. Environ. Sci.*, 10, <https://doi.org/10.3389/fenvs.2022.925979>, 2022.
- World Health Organization: WHO Global Air Quality Guidelines: Particulate Matter (PM<sub>2.5</sub> and PM<sub>10</sub>), Ozone, Nitrogen Dioxide, Sulfur Dioxide and Carbon Monoxide, World Health Organization, <https://www.who.int/publications/i/item/9789240034228> (last access: 8 February 2026), 2021.
- Wu, P.-C. and Huang, K.-F.: Tracing local sources and long-range transport of PM<sub>10</sub> in central Taiwan by using chemical characteristics and Pb isotope ratios, *Sci. Rep.*, 11, 7593, <https://doi.org/10.1038/s41598-021-87051-y>, 2021.
- Xiao, Q., Zheng, Y., Geng, G., Chen, C., Huang, X., Che, H., Zhang, X., He, K., and Zhang, Q.: Separating emission and meteorological contributions to long-term PM<sub>2.5</sub> trends over eastern China during 2000–2018, *Atmos. Chem. Phys.*, 21, 9475–9496, <https://doi.org/10.5194/acp-21-9475-2021>, 2021.
- Yin, H., Sun, Y., Liu, C., Zhang, L., Lu, X., Wang, W., Shan, C., Hu, Q., Tian, Y., Zhang, C., Su, W., Zhang, H., Palm, M., Notholt, J., and Liu, J.: FTIR time series of stratospheric NO<sub>2</sub> over Hefei, China, and comparisons with OMI and GEOS-Chem model data, *Opt. Express*, 27, A1225–A1240, <https://doi.org/10.1364/OE.27.0A1225>, 2019.
- Yin, H., Sun, Y., Liu, C., Lu, X., Smale, D., Blumenstock, T., Nagahama, T., Wang, W., Tian, Y., Hu, Q., Shan, C., Zhang, H., and Liu, J.: Ground-based FTIR observation of hydrogen chloride (HCl) over Hefei, China, and comparisons with GEOS-Chem model data and other ground-based FTIR stations data, *Opt. Express*, 28, 8041–8055, <https://doi.org/10.1364/OE.384377>, 2020.
- Yin, H., Sun, Y., Wang, W., Shan, C., Tian, Y., and Liu, C.: Ground-based high-resolution remote sensing of sulphur hexafluoride (SF<sub>6</sub>) over Hefei, China: characterization, optical misalignment, influence, and variability, *Opt. Express*, 29, 34051–34065, <https://doi.org/10.1364/OE.440193>, 2021a.
- Yin, H., Liu, C., Hu, Q., Liu, T., Wang, S., Gao, M., Xu, S., Zhang, C., and Su, W.: Opposite impact of emission reduction during the COVID-19 lockdown period on the surface concentrations of PM<sub>2.5</sub> and O<sub>3</sub> in Wuhan, China, *Environ. Pollut.*, 289, 117899, <https://doi.org/10.1016/j.envpol.2021.117899>, 2021b.
- Yin, H., Lu, X., Sun, Y., Li, K., Gao, M., Zheng, B., and Liu, C.: Unprecedented decline in summertime surface ozone over eastern China in 2020 comparably attributable to anthropogenic emission reductions and meteorology, *Environ. Res. Lett.*, 16, 124069, <https://doi.org/10.1088/1748-9326/ac3e22>, 2021c.
- Yin, H., Sun, Y., Notholt, J., Palm, M., Ye, C., and Liu, C.: Quantifying the drivers of surface ozone anomalies in the urban areas over the Qinghai-Tibet Plateau, *Atmos. Chem. Phys.*, 22, 14401–14419, <https://doi.org/10.5194/acp-22-14401-2022>, 2022a.
- Yin, H., Sun, Y., You, Y., Notholt, J., Palm, M., Wang, W., Shan, C., and Liu, C.: Using machine learning approach to reproduce the measured feature and understand the model-to-measurement discrepancy of atmospheric formaldehyde, *Sci. Total Environ.*, 851, 158271, <https://doi.org/10.1016/j.scitotenv.2022.158271>, 2022b.
- Zhang, H., Si, S., and Hsieh, C.-J.: GPU-acceleration for Large-scale Tree Boosting, *arXiv [preprint]*, <https://doi.org/10.48550/arXiv.1706.08359>, 2017.
- Zhang, Q., Meng, X., Shi, S., Kan, L., Chen, R., and Kan, H.: Overview of particulate air pollution and human health in China: Evidence, challenges, and opportunities, *The Innovation*, 3, 100312, <https://doi.org/10.1016/j.xinn.2022.100312>, 2022a.
- Zhang, R., Wang, G., Guo, S., Zamora, M. L., Ying, Q., Lin, Y., Wang, W., Hu, M., and Wang, Y.: Formation of urban fine particulate matter, *Chem. Rev.*, 115, 3803–3855, <https://doi.org/10.1021/acs.chemrev.5b00067>, 2015.
- Zhang, Y., Yu, S., Chen, X., Li, Z., Li, M., Song, Z., Liu, W., Li, P., Zhang, X., Lichtfouse, E., and Rosenfeld, D.: Local production, downward and regional transport aggravated surface ozone pollution during the historical orange-alert large-scale ozone episode in eastern China, *Environ. Chem. Lett.*, 20, 1577–1588, <https://doi.org/10.1007/s10311-022-01421-0>, 2022b.
- Zhao, R., Gu, X., Xue, B., Zhang, J., and Ren, W.: Short period PM<sub>2.5</sub> prediction based on multivariate linear regression model, *PLOS ONE*, 13, e0201011, <https://doi.org/10.1371/journal.pone.0201011>, 2018.
- Zhao, X. J., Zhao, P. S., Xu, J., Meng, W., Pu, W. W., Dong, F., He, D., and Shi, Q. F.: Analysis of a winter regional haze event and its formation mechanism in the North China Plain, *Atmos. Chem. Phys.*, 13, 5685–5696, <https://doi.org/10.5194/acp-13-5685-2013>, 2013.