

1 A Novel Method for Quantifying the Contribution of Regional Transport to PM_{2.5} in Beijing
2 (2013-2020): Combining Machine Learning with Concentration-Weighted Trajectory Analysis

3 Kang Hu¹, Hong Liao¹, Dantong Liu², Jianbing Jin¹, Lei Chen¹, Siyuan Li², Yangzhou Wu³,
4 Changhao Wu⁴, Shitong Zhao², Xiaotong Jiang⁵, Ping Tian^{6,7}, Kai Bi^{6,7}, Ye Wang⁸, Delong
5 Zhao^{6,7}

6 ¹Jiangsu Collaborative Innovation Center of Atmospheric Environment and Equipment
7 Technology, Jiangsu Key Laboratory of Atmospheric Environment Monitoring and Pollution
8 Control, Nanjing University of Information Science & Technology, Nanjing 210044, China.

9 ²Department of Atmospheric Sciences, School of Earth Sciences, Zhejiang University,
10 Hangzhou 310058, China.

11 ³Guangxi Key Laboratory of Environmental Pollution Control Theory and Technology, Guilin
12 University of Technology, Guilin 541004, China.

13 ⁴Institute of International Rivers and Eco-security, Yunnan University, Kunming 650091, China.

14 ⁵College of Biological and Environmental Engineering, Shandong University of Aeronautics,
15 Binzhou, 256600, China.

16 ⁶Beijing Key Laboratory of Cloud, Precipitation and Atmospheric Water Resources, Beijing
17 100089, China.

18 ⁷Field Experiment Base of Cloud and Precipitation Research in North China, China
19 Meteorological Administration, Beijing 100089, China.

20 ⁸Key Laboratory of Meteorological Disaster, Ministry of Education (KLME)/Joint
21 International Research Laboratory of Climate and Environment Change (ILCEC)/
22 Collaborative Innovation Center on Forecast and Evaluation of Meteorological Disasters (CIC-
23 FEMD), Nanjing University of Information Science and Technology, Nanjing 210044, China.

24

25

26 Corresponding author: Hong Liao (hongliao@nuist.edu.cn)

27

28

29 **Abstract**

30 Fine particulate matter (PM_{2.5}) is closely linked to human health, with its sources generally
31 divided into local emissions and regional transport. This study combined concentration-
32 weighted trajectory (CWT) analysis with the HYSPLIT trajectory ensemble to obtain hourly-
33 resolution pollutant source results. The Extreme Gradient Boosting (XGBoost) model was then
34 employed to simulate local emissions and ambient PM_{2.5} in Beijing from 2013 to 2020. The
35 results revealed that clean air masses influencing the Beijing area mainly originated from the
36 north and east regions, exhibiting a strong winter and weak summer pattern. Following the
37 implementation of the Air Pollution Prevention and Control Action Plan (Action Plan) by the
38 Chinese government in 2017, pollution in Beijing decreased significantly, with the most
39 substantial reduction in regional transport pollution events occurring in the west region during
40 summer. Regional transport pollution events were most frequent in spring, up to 1.8 times
41 higher than in winter. Pollutants mainly originated from the west and south regions, while
42 polluted air masses from the east showed the least reduction, and the proportion of pollution
43 sources from this region was gradually increasing. From 2013 to 2020, local emissions were
44 the main contributors to pollution events in Beijing. The Action Plan has more effectively
45 reduced pollution caused by regional transport, particularly during autumn and winter. This
46 finding underscores the importance of Beijing prioritizing local emission reduction while also
47 considering potential contributions from the east region to effectively mitigate pollution events.

48 **Keywords:** Fine particulate matter (PM_{2.5}); concentration-weighted trajectory (CWT);
49 XGBoost model; regional transport

50

51 **1. Introduction**

52 Ambient fine particulate matter (PM_{2.5}, with particle aerodynamic diameter $\leq 2.5 \mu\text{m}$) is
53 influenced by both natural sources, such as dust, volcanic eruptions, tsunamis, and forest fires,
54 and anthropogenic emissions, including fuel combustion, transportation, and industrial
55 production. Anthropogenic emissions dominate the long-term trend of air pollution (Zhang et
56 al., 2019; Cheng et al., 2019). Numerous epidemiological studies have found that PM_{2.5} can
57 significantly damage human health by exacerbating respiratory and cardiovascular diseases
58 (Bartell et al., 2013; Brauer et al., 2012; Pascal et al., 2014), and also has an impact on weather
59 and climate change (Wang et al., 2014; Smith et al., 2020; Kalisoras et al., 2023). China's rapid
60 and energy-intensive development over the past several decades has led to severe air pollution
61 and negative public health impacts (Huang et al., 2014; Geng et al., 2021). Consequently,
62 controlling pollution and reducing PM_{2.5} concentrations became an urgent issue in China. While
63 meteorological variations caused about 16% of the ambient PM_{2.5} decline during 2013-2017
64 (Zhang et al., 2019), the uncertainty in reducing PM_{2.5} through meteorological conditions is
65 substantial, and the magnitude of the decrease is not dominated by human actions. Thus, the
66 primary means of controlling PM_{2.5} relies on reducing anthropogenic emissions. To address this
67 issue, the Chinese government implemented the Air Pollution Prevention and Control Action
68 Plan (denoted "Action Plan") from 2013 to 2017 and the Blue Sky Protection Campaign from
69 2018 to 2020, which effectively controlled anthropogenic emissions and reduced ambient PM_{2.5}
70 concentrations (Zhang et al., 2019; Du et al., 2022).

71 The concentration of PM_{2.5} can be attributed to local emissions and regional transport. Several
72 methods, such as the HYSPLIT model (Draxler and Rolph, 2010), can be used to distinguish
73 pollutant sources. Wu et al. used the HYSPLIT model to simulate the 24-hour backward
74 trajectory in Zhoushan (Wu et al., 2021), and identified continental air masses that spent more
75 than 5% of the previous 24 hours over the continent region, while the remaining air masses
76 were identified as oceanic-influenced air masses. Ding et al. employed a backward trajectory
77 ensemble to analyze the sources of air masses in Beijing during the study period (Ding et al.,
78 2019), finding that air masses with high concentrations of black carbon (BC) mass mainly came
79 from the south and southeast regions. Cluster analysis on backward trajectories can be used to
80 obtain the main direction of aerosols over a period of time, allowing for the analysis and
81 determination of dominant air mass directions. For instance, Li et al. divided the sources of air
82 masses in the Wuhan area from October to November 2019 into short transport distance,
83 northbound air masses, and regional transport from the northeast and some coastal areas (Li et
84 al., 2022).

85 The HYSPLIT model results are mainly used to view air mass trajectories, making it difficult
86 to directly determine the sources of pollutants. Potential source contribution function (PSCF)
87 and concentration-weighted trajectory (CWT) analyses based on backward trajectories can be
88 used to identify the sources of pollutants through conditional probability results. Hu et al. used
89 weighted PSCF to analyze the sources of air masses with different levels of pollution in Beijing
90 and found that polluted air masses from the southwest were an important source of high-level
91 advections during the study period, while light pollution was often accompanied by the regional
92 transport originating from the northeast region (Hu et al., 2020). Wu et al. used CWT to analyze
93 the sources of pollution in Zhoushan and found that pollutants in Zhoushan are influenced by

94 both local emissions and regional transport. There are no obvious high pollution areas, while in
95 other seasons, PM_{2.5} mainly originates from southern Jiangsu and Shanghai (Wu et al., 2024).
96 However, these studies relied on standard HYSPLIT trajectory results, which have lower
97 temporal resolution, limiting the accuracy of pollutant source identification.

98 The Lagrangian air pollution dispersion model, Numerical Atmospheric-dispersion Modelling
99 Environment (NAME) (Jones et al., 2007) can determine the source of polluted air masses by
100 simulating particulate concentrations within each grid point using Monte Carlo methods,
101 followed by 3-D trajectories of plume basins. Liu et al. used the NAME model to study the
102 sources of air masses in Beijing during the winter of 2019 and divided them into local emissions
103 and regional transport to analyze the convective mixing process of BC under the influence of
104 local emissions (Liu et al., 2020). However, due to limitations in computing resources, the
105 NAME model is difficult to use for obtaining long-term emission source analysis results.

106 Multiple methods can be used to predict PM_{2.5} concentrations, such as statistical models (e.g.,
107 linear mixed-effect models and generalized additive models) (Fang et al., 2016; Ma et al., 2016),
108 chemical transport model (CTM)-based algorithms (Geng et al., 2015; Kong et al., 2021),
109 physical models (Lin et al., 2018), and recently emerging machine learning models, including
110 Extreme Gradient Boosting (XGBoost) and Random Forest (Liang et al., 2020; Wei et al., 2021;
111 Xiao et al., 2018; Xue et al., 2019; Huang et al., 2021). Geng et al. used satellite observations
112 of aerosol optical depth (AOD) and meteorological data combined with the XGBoost model to
113 explore the long-term variations of PM_{2.5} caused by changes in meteorological conditions from
114 2000 to 2018 (Geng et al., 2021). Kleine Deters et al. demonstrated the relevance of statistical
115 models based on machine learning for predicting PM_{2.5} concentrations from meteorological
116 data (Kleine Deters et al., 2017). This method of predicting aerosol concentrations using only
117 meteorological data has been widely used (Asadollahfardi et al., 2016; Zeng et al., 2021). For
118 instance, Grange et al. used meteorological data, synoptic scale weather patterns, and time
119 variables to explain daily PM₁₀ concentrations in Switzerland (Grange et al., 2018). In summary,
120 machine learning models have achieved high accuracy in estimating and predicting PM_{2.5}
121 concentrations and have high use value, and the rise of machine learning methods has also
122 provided feasibility for quantifying the contribution of regionally transported air masses.

123 In this study, we combined CWT analysis with the HYSPLIT trajectory ensemble to obtain
124 hourly-resolution PM_{2.5} source results and used this approach to distinguish between local
125 emissions and regional transport. Solved the problems of traditional CWT methods being
126 unable to obtain hourly time accuracy and models such as NAME consuming a large number
127 of computational resources. Predictive XGBoost models were developed for Beijing using
128 meteorological data and time variables to explain PM_{2.5} concentrations. By training the
129 XGBoost model with PM_{2.5} dominated by local emissions, which are separately distinguished
130 by CWT, and generalizing the findings to all study periods, the concentration of locally emitted
131 PM_{2.5} (local) can be obtained. Similarly, ambient observed PM_{2.5} (ambient) can be determined
132 by training the XGBoost model with ambient PM_{2.5} data. The contribution of regional transport
133 to PM_{2.5} in Beijing can be quantified by comparing the ambient and local PM_{2.5} concentrations.

134

135 2. Materials and methods

136 2.1 Site and instrumentation

137 The PM_{2.5} data (Fig. 1a) were obtained from in situ air quality monitoring conducted by the
138 China National Environmental Monitoring Center from 2013 to 2020. The monitoring station
139 is located in Haidian Wanliu (39.96°N, 116.29°E), situated in the central urban area of Beijing.
140 Meteorological data, including temperature, relative humidity, pressure, precipitation, wind
141 speed, and planetary boundary layer height (PBLH), were sourced from the European Centre
142 for Medium-Range Weather Forecasts (ECMWF) ERA5 hourly reanalysis dataset
143 (<https://cds.climate.copernicus.eu/datasets>). In this study, a year was divided into four quarters:
144 Spring (March, April, and May), Summer (June, July, and August), Autumn (September,
145 October, and November), and Winter (December, January, and February).

146

147 2.2 Air mass source

148 The air mass trajectory data were obtained from the 1°×1° horizontal and vertical wind fields
149 of the Global Data Assimilation System (GDAS) reanalysis products
150 (<ftp://arlftp.arlhq.noaa.gov/pub/archives/gdas1>), which are available every 3 hours. The
151 HYSPLIT trajectory ensemble was used to generate 27 equally probable 24-hour backward air
152 mass trajectories for the target point (39.96°N, 116.29°E, 250 m a.s.l.) in every hour by using
153 PySplit (Cross, 2015). Given the equal probability of air masses being transported to the target
154 point for each trajectory in the HYSPLIT trajectory ensemble, a conditional probability CWT
155 analysis was applied to determine the hourly source area of pollution.

156 In the CWT analysis method, each grid point is assigned a weight, and the contribution of each
157 grid point to the pollutant concentration at the target site is calculated using the air mass
158 residence time and pollutant concentration (Hopke et al., 1993; Polissar et al., 1999; Xu and
159 Akhtar, 2010) (equation 1). The grid point resolution was set to 0.25°×0.25° for this study. In
160 equations 1, C_{ij} is the average weighted concentration at grid point (i, j) , l is the trajectory
161 index, M represents the total number of trajectories, C_l is the PM_{2.5} concentration
162 corresponding to the target site, and τ_{ijl} is the residence time of trajectory l passing through
163 the grid point. In calculation, the number of trajectories falling on each grid point is used instead
164 of the residence time.

$$165 \quad C_{ij} = \frac{\sum_{l=1}^M C_l \times \tau_{ijl}}{\sum_{l=1}^M \tau_{ijl}} \quad (1)$$

166 To reduce the effect of small values of n_{ij} , the CWT values were multiplied by an arbitrary
167 weight function $W(n_{i,j})$ to better reflect the uncertainty in the values for these grids (equation
168 2).

$$169 \quad W(n_{i,j}) = \begin{cases} 1.00, & 3n_{ave} < n_{ij} \\ 0.70, & 1.5n_{ave} < n_{ij} \leq 3n_{ave} \\ 0.4, & n_{ave} < n_{ij} \leq 1.5n_{ave} \\ 0.17, & n_{ij} \leq n_{ave} \end{cases} \quad (2)$$

170 where n_{ij} represents the number of trajectories that fall within the grid point, and n_{ave}
171 represents the average number of trajectories passing through each grid point.

172 The potential source contribution to PM_{2.5} at the target site was investigated by categorizing the
173 backward air masses into five different source regions centered around Beijing: local (which is
174 a region around central Beijing, 115.3~117.5°E, 39.4~41°N); north region (the northern plateau
175 at 108~117.5°E, 41~43°N); west region (the western plateau at 108~115.3°E, 34~41°N); south
176 region (the southern plain at 115.3~120°E, 34~39.4°N); and east region (the eastern plain at
177 117.5~120°E, 39.4~43°N). The concentration is integrated over each grid point in each
178 segregated region obtained from the CWT analysis, and the contributions of each air mass
179 fraction are obtained. The region with the highest contribution is used to determine the
180 dominant source of air masses in Beijing at each time, classifying the overall air mass sources
181 into local emissions (Fig. 1g) and regional transport (Fig. 1h). It is important to note that local
182 emission periods were also influenced by persistent regional transport, and vice versa.

183

184 2.3 Deriving the long-term local emission and ambient PM_{2.5}

185 An XGBoost model is employed to derive the local and ambient PM_{2.5} results. The
186 hyperparameters used in the model for local (ambient) conditions include a maximum number
187 of boosting iterations of 6067 (13421), a learning rate of 0.1, a maximum tree depth of 7 (11),
188 a minimum sum of instance weight needed in a child of 5 (3), a subsampling ratio of 0.8 (0.6)
189 for training instances, and a subsampling ratio of 0.8 for columns when constructing each tree.
190 The input parameters for the XGBoost model comprise meteorological variables (temperature,
191 relative humidity, wind speed, surface pressure, and precipitation) and temporal parameters
192 (year, month, day of the week, and day of the year), as referenced from Xu et al. (Xu et al.,
193 2023). Additionally, PBLH, which has been shown to significantly impact pollutant
194 concentrations in previous observational (Su et al., 2018; Miao and Liu, 2019; Miao et al., 2019)
195 and machine learning studies (Xiao et al., 2021; Li et al., 2017b; Shen et al., 2018), was included
196 as an input parameter. Based on the XGBoost learning results, the most sensitive parameters
197 for both local and ambient PM_{2.5} are RH, wind field, surface pressure and PBLH (Fig. S1). For
198 the machine learning process, data from 2013 to 2019 were used for training the XGBoost
199 models, while data from 2020 were used for model validation. Note that the 2020 analysis
200 results may contain some uncertainties due to the impact of COVID-19.

201 The relatively small proportion of high-concentration PM_{2.5} can lead to underestimation of
202 high-concentration events in the model results (Wei et al., 2020). To address this issue, a high
203 PM_{2.5} indicator was defined as a daily average PM_{2.5} concentration exceeding the monthly
204 average plus twice the standard deviation. In this study, original high PM_{2.5} indicators accounted
205 for 6% of the data points during the period dominated by local and ambient PM_{2.5}. To balance
206 the proportion of high-concentration PM_{2.5} in the entire database, the Synthetic Minority Over-
207 sampling Technique (SMOTE) (Torgo, 2011) was applied during data preprocessing. SMOTE
208 artificially generates new synthetic samples along the line between high-concentration data
209 points and their selected nearest neighbors, effectively oversampling the high-concentration
210 data. As a result, the proportion of high PM_{2.5} indicators increased to 21% and 22% for local
211 and ambient PM_{2.5}, respectively.

212 Hyperparameter optimization and performance evaluation of the model were conducted using
213 fivefold cross-validation (CV), while early stopping with a patience of 10 rounds was employed
214 to prevent overfitting. (Akritidis et al., 2021; Zhang et al., 2020). In this approach, 20% of the
215 data is randomly selected for model testing, while the remaining 80% is used for training. This
216 process is repeated five times, ensuring that each record is used once as testing data. The
217 coefficient of determination (r^2) was employed to assess the correlation between the XGBoost
218 model predictions and observed values, while the root mean square error (RMSE) was used as
219 a performance evaluation statistic. After obtaining the relation between the input parameters
220 and $PM_{2.5}$, we are able to derive the hourly local and ambient $PM_{2.5}$ once all long-term input
221 parameters (Fig. S3).

222 3 Results and discussion

223 3.1 Evaluation of the XGBoost $PM_{2.5}$ prediction model

224 During the model validation process, the XGBoost model results for ambient $PM_{2.5}$ (Fig. 2a2)
225 demonstrated an r^2 of 0.74 and an RMSE of $20 \mu\text{g m}^{-3}$ when compared to observations. The
226 XGBoost model results for local $PM_{2.5}$ exhibited an r^2 of 0.78 and an RMSE of $21 \mu\text{g m}^{-3}$. An
227 analysis of the $PM_{2.5}$ frequency distribution in Beijing revealed an agreement between the
228 XGBoost model results and observations for both ambient and local $PM_{2.5}$. The error between
229 XGBoost learning results and actual observed $PM_{2.5}$ values is mainly concentrated in the low
230 concentration stage. This may be attributed to the significant reduction in human activities
231 during the COVID-19 lockdown periods, which led to a decrease in actual $PM_{2.5}$ levels, making
232 it challenging for XGBoost to learn (Fig. 2b1 and b2). As illustrated in Fig. S2, local and
233 ambient $PM_{2.5}$ in Beijing display a distinct seasonal variation, with higher values in winter and
234 lower values in summer. However, the transport of clean air masses from the north diminishes
235 the seasonal variation characteristics of ambient $PM_{2.5}$ in Beijing, making winter pollution less
236 prominent compared to other seasons.

237 Fig. S3 reveals that ambient pollution events ($PM_{2.5} > 75 \mu\text{g m}^{-3}$) in Beijing are primarily
238 influenced by air masses originating from the south and west, particularly under the control of
239 westward air masses. Numerous studies have indicated that air masses originating from the
240 western region significantly contribute to regional pollution events in Beijing (Streets et al.,
241 2007; Tian et al., 2019; Liu et al., 2020). With the exception of December (Fig. 3b1), westward
242 air masses often bring higher monthly average $PM_{2.5}$ to Beijing. Air masses originating from
243 the south region can also transport more pollutants to Beijing (Fig. S3). However, unlike the
244 high-frequency polluted air masses from the west, southward air masses are associated with
245 higher $PM_{2.5}$ concentrations, particularly during autumn and winter (Fig. 3c1). This
246 phenomenon can be attributed to the higher pollution levels in Hebei and Shandong provinces
247 compared to Beijing during these seasons, as verified by AOD observations from Moderate
248 Resolution Imaging Spectroradiometer (MODIS) on the Aqua satellites over Eastern China
249 (Zhang and Reid, 2010; Hu et al., 2018) (Fig. S4). Notably, in contrast to westward transport,
250 air masses from the south region in February predominantly exhibited a cleaning effect on
251 Beijing, even before 2017 (Fig. S3b). This can be explained by the occurrence of these transport
252 processes during or shortly after the Spring Festival, a period characterized by extremely low

253 anthropogenic emissions, resulting in lower ambient PM_{2.5} compared to local emissions in the
254 megacity of Beijing. Following the implementation of the Action Plan, the polluted air masses
255 from the south region transitioned from carrying higher PM_{2.5} to levels close to local emission
256 concentrations in Beijing, leading to a more equal contribution to pollution and clean events in
257 the area (Fig. S5c1).

258 3.2 Impact of clean air masses from transported regions on PM_{2.5} in Beijing

259 In this study, clean air masses are defined as those associated with ambient PM_{2.5} in the Beijing
260 area that are lower than the concentrations resulting from local emissions, as illustrated below
261 the dashed line in Fig. 3a1-d1. This study reveals that clean air masses predominantly originate
262 from the east and north regions during the period 2013-2020, which is consistent with previous
263 studies (Zhang et al., 2018; Hu et al., 2020). Clean air masses from different directions exhibit
264 similar seasonal variations in their ability to reduce locally emitted pollution in Beijing, with a
265 strong reduction effect in winter and a weaker effect in summer (Fig. 3a2-d2). This
266 phenomenon is closely related to the seasonal variations in pollutant emissions. Due to the
267 combined influence of increased residential emissions from heating activities and
268 meteorological conditions in Beijing during autumn and winter, local PM_{2.5} in Beijing presents
269 higher concentrations. Consequently, the influx of clean air masses results in a more
270 pronounced reduction in PM_{2.5} during these seasons. The weaker attenuation effect of PM_{2.5}
271 transported from the south region during December and January can be attributed to the high-
272 frequency and high-concentration pollution contributions from air masses originating in this
273 region during this period.

274 Due to a significant reduction in anthropogenic emissions after 2017, the attenuation of PM_{2.5}
275 concentrations by clean air masses from all directions was significantly lower than before 2017
276 (Fig. S6a2-d2). Compared to the period prior to 2017, the mean attenuation of PM_{2.5}
277 concentrations in Beijing decreased by 3, 10, 3, and 7 $\mu\text{g m}^{-3}$ ($p < 0.01$) for air masses
278 originating from the north, west, south, and east regions, respectively.

279 3.3 Variations in Beijing PM_{2.5} concentrations under transport-induced pollution events

280 Transport-induced pollution events in Beijing are defined as the occurrence of ambient PM_{2.5}
281 exceeding both local PM_{2.5} and the light pollution standard ($75 \mu\text{g m}^{-3}$). Fig. 4a1-d1 demonstrate
282 that the monthly variation of PM_{2.5} in Beijing generally follows a unimodal pattern, with higher
283 values in winter and lower values in summer, except when under the influence of eastern air
284 mass transport. This phenomenon is closely related to the seasonal variations in anthropogenic
285 emissions in China and the characteristics of climate change (Renhe et al., 2014; Li et al., 2017a;
286 Zhang et al., 2015). The overall PM_{2.5} in Beijing under the influence of eastward pollution air
287 masses exhibits a bimodal distribution, with frequent high-concentration pollution events
288 occurring in January and October. Even after the effective control of anthropogenic emissions
289 in 2017, a second peak of high-concentration pollution persists in October (Fig. 4d2). Fig. 4a2-
290 d2 illustrate the effectiveness of the Action Plan in controlling pollutant concentrations in the
291 Beijing area. Since 2017, PM_{2.5} in Beijing has been significantly lower than the values observed
292 before 2017 during transport-induced pollution events. Moreover, during January and from

293 June to September, there were periods when the regional transport of polluted air masses from
294 a fixed direction did not contribute to pollution events in Beijing.

295 An analysis of the proportion of transport-induced pollution events from different regions to
296 Beijing (Fig. 5) shows that after the implementation of the Action Plan in 2017, the number of
297 pollution events dominated by regional transport decreased significantly. From spring to winter,
298 the largest decrease in transport-induced pollution events occurred in the north, west, west and
299 south regions in each season, with the lowest decrease occurring in the east region during winter.

300 The temporal variation in the number of transport-induced pollution events from different
301 regions (Fig. S7) revealed that air masses transported from the west region contributed to the
302 most frequent pollution events in each season except summer. The highest number of events
303 occurred in spring 2016 (322), autumn 2016 (375), and winter 2017 (308). Summer transport-
304 induced pollution events were mainly influenced by polluted air masses transported from the
305 south, with a gradual decrease in the number of events over the years. Although pollution events
306 in Beijing primarily occur in autumn and winter, this study found that after 2017, the season
307 when Beijing was most affected by transport-induced pollution events was spring, contributing
308 a total of 685 pollution events, while autumn and winter contributed 266 and 392 events,
309 respectively. The impact of polluted air masses on summer transport was minimal, with only
310 215 occurrences.

311 Fig. 5a shows that in spring, transport-induced pollution events in Beijing were mainly
312 dominated by polluted air masses transported from the west and south. The highest proportion
313 of regional transport events from the west occurred in 2016, reaching 68%, while the highest
314 proportion of southward transport-induced pollution events occurred in 2017 (with the
315 exception of 2020, which may have been influenced by the COVID-19 pandemic). The
316 increased frequency of pollution air masses transported from the south after 2017 can be
317 attributed to the effective control of anthropogenic emissions, resulting in a decrease in PM_{2.5}
318 transported from various regions, especially from westward sources (Fig. S7a). The decrease
319 in the proportion of pollution events transported from the west, which originally accounted for
320 a large proportion, led to an increase in the contribution of remaining incoming air masses to
321 Beijing.

322 Before 2017, transport-induced pollution events in Beijing during summer were mainly
323 affected by polluted air masses from the south region. Even in 2015, when the proportion of
324 transport-induced pollution events from south region was lowest during the entire period, it still
325 accounted for 50% of the total number of transport-induced pollution events that year. However,
326 after the implementation of the Action Plan, the proportion of transport-induced pollution
327 events from the south region gradually decreased to 38%. In 2020, this proportion further
328 declined to 25%, but this may have been affected by the COVID-19 pandemic. Meanwhile,
329 pollution air masses originating from the east increasingly dominated the occurrence of
330 pollution events in Beijing.

331 Transport-induced pollution events in Beijing mainly originated from the west and had the
332 highest contribution proportion in autumn before 2019 (except for 2013, when the contribution

333 proportion was 34%, second only to southward air masses at 35%). After 2019, the contribution
334 of eastward air masses became dominant in autumn. In winter, polluted air masses from the
335 west were the main source of transport-induced pollution events. Overall, as the Action Plan
336 gradually improved, the transport-induced pollution from the east did not decrease significantly
337 compared to other air mass sources. This may be because the eastward air masses are mostly
338 clean. However, as the concentration of polluted air masses from other sources decreases, the
339 potential impact of eastward air masses on Beijing's transport-induced pollution events
340 increases. This finding may prompt Beijing to prioritize emission reduction in the east region
341 when implementing future joint prevention and control measures.

342 **4 Conclusion**

343 This study combined a machine learning method and Concentration-Weighted Trajectory
344 (CWT) analysis to derive local emissions and ambient observed PM_{2.5} in Beijing from 2013 to
345 2020, thus the contribution of regional transport to PM_{2.5} in Beijing can be quantified. The
346 impact of clean air masses (defined as those with ambient PM_{2.5} concentrations lower than local
347 emissions) mainly originated from the east and north regions. These clean air masses from
348 different directions exhibited similar seasonal variations in their ability to reduce ambient
349 pollution in Beijing, with a stronger reduction effect in winter and a weaker reduction effect in
350 summer.

351 Except for the regional transport from the east region, the seasonal variation of PM_{2.5} in Beijing
352 under the influence of transport-induced pollution events (ambient PM_{2.5} exceeding both local
353 PM_{2.5} and 75 $\mu\text{g m}^{-3}$) shows a general trend of high concentrations in winter and low
354 concentrations in summer. The main reason for this phenomenon is related to the seasonal
355 emissions of pollutants in China and the characteristics of climate change. Before 2019, the
356 west region was the primary source of pollution events during autumn and winter. However,
357 starting from 2019, the east region became the main contributor of polluted air masses in
358 autumn. Additionally, among all regions, the east region exhibited the smallest decrease in
359 transport-induced pollution events after 2017.

360 From 2013 to 2020, local emissions were the main contributors to pollution events in Beijing.
361 However, the Air Pollution Prevention and Control Action Plan, implemented by the Chinese
362 government in 2017, more effectively mitigated pollutants caused by regional transport
363 compared to local emissions, particularly during autumn and winter. This finding suggests that
364 Beijing should prioritize reducing local emissions while also accounting for potential
365 contributions from the east region in its future pollution prevention and control strategies.

366

367 **Code and data availability**

368 The codes used in this study are archived on Zenodo: the machine learning code at
369 <https://doi.org/10.5281/zenodo.14677125>, the CWT code at
370 <https://doi.org/10.5281/zenodo.13994400>, ECMWF data at

371 <https://doi.org/10.5281/zenodo.14353871>, GDAS data at
372 <https://doi.org/10.5281/zenodo.14347277>, HySplit Trajectory Ensemble at
373 <https://doi.org/10.5281/zenodo.14375567>, and PySPLIT at
374 <https://doi.org/10.5281/zenodo.14354765>. The meteorology and PM_{2.5} data used in this study
375 can be accessed at <https://dx.doi.org/10.17632/bhfktx3kz8.2>.

376 **Author contribution**

377 Kang Hu, Hong Liao and Dantong Liu designed and carried out the experiments. Kang Hu
378 wrote the code and final paper with contributions from all other authors. Hong Liao, Dantong
379 Liu, Lei Chen and Jianbing Jin reviewed and edited the paper.

380

381 **Competing interests**

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

383

384 **Acknowledgements**

385 This research was supported by the China Postdoctoral Science Foundation (2023M741773),
386 Postdoctoral Fellowship Program of CPSF (GZC20231150), National Natural Science
387 Foundation of China (42021004, 42405192).

388

389 **Reference**

390 Akritidis, D., Zanis, P., Georgoulas, A. K., Papakosta, E., Tzoumaka, P., and Kelessis, A.:
391 Implications of COVID-19 restriction measures in urban air quality of Thessaloniki, Greece: A
392 machine learning approach, *Atmosphere*, 12, 1500, 2021.

393 Asadollahfardi, G., Madinejad, M., Aria, S. H., and Motamadi, V.: Predicting Particulate Matter
394 (PM_{2.5}) Concentrations in the Air of Shahr-e Ray City, Iran, by Using an Artificial Neural
395 Network, *Environmental Quality Management*, 25, 71-83, 2016.

396 Bartell, S. M., Longhurst, J., Tjoa, T., Sioutas, C., and Delfino, R. J.: Particulate air pollution,
397 ambulatory heart rate variability, and cardiac arrhythmia in retirement community residents
398 with coronary artery disease, *Environmental health perspectives*, 121, 1135-1141, 2013.

399 Brauer, M., Amann, M., Burnett, R. T., Cohen, A., Dentener, F., Ezzati, M., Henderson, S. B.,
400 Krzyzanowski, M., Martin, R. V., and Van Dingenen, R.: Exposure assessment for estimation
401 of the global burden of disease attributable to outdoor air pollution, *Environmental science &*
402 *technology*, 46, 652-660, 2012.

403 Cheng, N., Cheng, B., Li, S., and Ning, T.: Effects of meteorology and emission reduction
404 measures on air pollution in Beijing during heating seasons, *Atmospheric Pollution Research*,
405 10, 971-979, 2019.

406 Cross, M.: PySPLIT: a Package for the Generation, Analysis, and Visualization of HYSPLIT
407 Air Parcel Trajectories, *SciPy*, 133-137,

408 Ding, S., Zhao, D., He, C., Huang, M., He, H., Tian, P., Liu, Q., Bi, K., Yu, C., and Pitt, J.:
409 Observed interactions between black carbon and hydrometeor during wet scavenging in mixed-
410 phase clouds, *Geophysical Research Letters*, 46, 8453-8463, 2019.

411 Draxler, R. and Rolph, G.: HYSPLIT (HYbrid Single-Particle Lagrangian Integrated Trajectory)
412 model access via NOAA ARL READY website (<http://ready.arl.noaa.gov/HYSPLIT.php>).
413 NOAA Air Resources Laboratory, Silver Spring, MD, 25, 2010.

414 Du, H., Li, J., Wang, Z., Chen, X., Yang, W., Sun, Y., Xin, J., Pan, X., Wang, W., and Ye, Q.:
415 Assessment of the effect of meteorological and emission variations on winter PM_{2.5} over the
416 North China Plain in the three-year action plan against air pollution in 2018–2020, *Atmospheric*
417 *Research*, 280, 106395, 2022.

418 Fang, X., Zou, B., Liu, X., Sternberg, T., and Zhai, L.: Satellite-based ground PM_{2.5} estimation
419 using timely structure adaptive modeling, *Remote Sensing of Environment*, 186, 152-163, 2016.

420 Geng, G., Zhang, Q., Martin, R. V., van Donkelaar, A., Huo, H., Che, H., Lin, J., and He, K.:
421 Estimating long-term PM_{2.5} concentrations in China using satellite-based aerosol optical
422 depth and a chemical transport model, *Remote sensing of Environment*, 166, 262-270, 2015.

423 Geng, G., Xiao, Q., Liu, S., Liu, X., Cheng, J., Zheng, Y., Xue, T., Tong, D., Zheng, B., and
424 Peng, Y.: Tracking air pollution in China: near real-time PM_{2.5} retrievals from multisource
425 data fusion, *Environmental Science & Technology*, 55, 12106-12115, 2021.

426 Grange, S. K., Carslaw, D. C., Lewis, A. C., Boleti, E., and Hueglin, C.: Random forest
427 meteorological normalisation models for Swiss PM₁₀ trend analysis, *Atmospheric Chemistry*
428 *and Physics*, 18, 6223-6239, 2018.

429 Hopke, P. K., Gao, N., and Cheng, M.-D.: Combining chemical and meteorological data to infer
430 source areas of airborne pollutants, *Chemometrics and Intelligent Laboratory Systems*, 19, 187-
431 199, 1993.

432 Hu, K., Kumar, K. R., Kang, N., Boiyo, R., and Wu, J.: Spatiotemporal characteristics of
433 aerosols and their trends over mainland China with the recent Collection 6 MODIS and OMI
434 satellite datasets, *Environmental Science and Pollution Research*, 25, 6909-6927, 2018.

435 Hu, K., Zhao, D., Liu, D., Ding, S., Tian, P., Yu, C., Zhou, W., Huang, M., and Ding, D.:
436 Estimating radiative impacts of black carbon associated with mixing state in the lower
437 atmosphere over the northern North China Plain, *Chemosphere*, 252, 126455, 2020.

438 Huang, C., Hu, J., Xue, T., Xu, H., and Wang, M.: High-resolution spatiotemporal modeling for
439 ambient PM_{2.5} exposure assessment in China from 2013 to 2019, *Environmental Science &*
440 *Technology*, 55, 2152-2162, 2021.

441 Huang, R.-J., Zhang, Y., Bozzetti, C., Ho, K.-F., Cao, J.-J., Han, Y., Daellenbach, K. R., Slowik,
442 J. G., Platt, S. M., and Canonaco, F.: High secondary aerosol contribution to particulate
443 pollution during haze events in China, *Nature*, 514, 218-222, 2014.

444 Jones, A., Thomson, D., Hort, M., and Devenish, B.: The UK Met Office's next-generation
445 atmospheric dispersion model, NAME III, in: *Air pollution modeling and its application XVII*,
446 Springer, 580-589, 2007.

447 Kalisoras, A., Georgoulas, A. K., Akritidis, D., Allen, R. J., Naik, V., Kuo, C., Szopa, S., Nabat,
448 P., Olivie, D., and Van Noije, T.: Decomposing the Effective Radiative Forcing of
449 anthropogenic aerosols based on CMIP6 Earth System Models, *Atmospheric Chemistry &*
450 *Physics Discussions*, 2023.

451 Kleine Deters, J., Zalakeviciute, R., Gonzalez, M., and Rybarczyk, Y.: Modeling PM_{2.5} urban
452 pollution using machine learning and selected meteorological parameters, *Journal of Electrical*
453 *and Computer Engineering*, 2017, 5106045, 2017.

454 Kong, L., Tang, X., Zhu, J., Wang, Z., Li, J., Wu, H., Wu, Q., Chen, H., Zhu, L., and Wang, W.:
455 A 6-year-long (2013–2018) high-resolution air quality reanalysis dataset in China based on the
456 assimilation of surface observations from CNEMC, *Earth System Science Data*, 13, 529-570,
457 2021.

458 Li, M., Zhang, Q., Kurokawa, J.-i., Woo, J.-H., He, K., Lu, Z., Ohara, T., Song, Y., Streets, D.
459 G., and Carmichael, G. R.: MIX: a mosaic Asian anthropogenic emission inventory under the
460 international collaboration framework of the MICS-Asia and HTAP, *Atmospheric Chemistry*
461 *and Physics*, 17, 935-963, 2017a.

462 Li, S., Liu, D., Kong, S., Wu, Y., Hu, K., Zheng, H., Cheng, Y., Zheng, S., Jiang, X., and Ding,
463 S.: Evolution of source attributed organic aerosols and gases in a megacity of central China,
464 *Atmospheric Chemistry and Physics Discussions*, 2022, 1-19, 2022.

465 Li, T., Shen, H., Yuan, Q., Zhang, X., and Zhang, L.: Estimating ground-level PM_{2.5} by fusing
466 satellite and station observations: a geo-intelligent deep learning approach, *Geophysical*
467 *Research Letters*, 44, 11,985-911,993, 2017b.

468 Liang, F., Xiao, Q., Huang, K., Yang, X., Liu, F., Li, J., Lu, X., Liu, Y., and Gu, D.: The 17-y
469 spatiotemporal trend of PM_{2.5} and its mortality burden in China, *Proceedings of the National*
470 *Academy of Sciences*, 117, 25601-25608, 2020.

471 Lin, C., Liu, G., Lau, A. K. H., Li, Y., Li, C., Fung, J. C. H., and Lao, X. Q.: High-resolution
472 satellite remote sensing of provincial PM_{2.5} trends in China from 2001 to 2015, *Atmospheric*
473 *environment*, 180, 110-116, 2018.

474 Liu, D., Hu, K., Zhao, D., Ding, S., Wu, Y., Zhou, C., Yu, C., Tian, P., Liu, Q., and Bi, K.:
475 Efficient vertical transport of black carbon in the planetary boundary layer, *Geophysical*
476 *Research Letters*, 47, e2020GL088858, 2020.

477 Ma, Z., Hu, X., Sayer, A. M., Levy, R., Zhang, Q., Xue, Y., Tong, S., Bi, J., Huang, L., and Liu,
478 Y.: Satellite-based spatiotemporal trends in PM_{2.5} concentrations: China, 2004–2013,
479 *Environmental health perspectives*, 124, 184-192, 2016.

480 Miao, Y. and Liu, S.: Linkages between aerosol pollution and planetary boundary layer structure
481 in China, *Science of the Total Environment*, 650, 288-296, 2019.

482 Miao, Y., Li, J., Miao, S., Che, H., Wang, Y., Zhang, X., Zhu, R., and Liu, S.: Interaction
483 between planetary boundary layer and PM_{2.5} pollution in megacities in China: a Review,
484 *Current Pollution Reports*, 5, 261-271, 2019.

485 Pascal, M., Falq, G., Wagner, V., Chatignoux, E., Corso, M., Blanchard, M., Host, S., Pascal,
486 L., and Larrieu, S.: Short-term impacts of particulate matter (PM₁₀, PM_{10-2.5}, PM_{2.5}) on
487 mortality in nine French cities, *Atmospheric Environment*, 95, 175-184, 2014.

488 Polissar, A., Hopke, P., Paatero, P., Kaufmann, Y., Hall, D., Bodhaine, B., Dutton, E., and Harris,
489 J.: The aerosol at Barrow, Alaska: long-term trends and source locations, *Atmospheric*
490 *Environment*, 33, 2441-2458, 1999.

491 Renhe, Z., Li, Q., and Zhang, R.: Meteorological conditions for the persistent severe fog and
492 haze event over eastern China in January 2013, *Science China Earth Sciences*, 57, 26-35, 2014.

493 Shen, H., Li, T., Yuan, Q., and Zhang, L.: Estimating regional ground-level PM_{2.5} directly
494 from satellite top-of-atmosphere reflectance using deep belief networks, *Journal of Geophysical*
495 *Research: Atmospheres*, 123, 875-813,886, 2018.

496 Smith, C. J., Kramer, R. J., Myhre, G., Alterskjær, K., Collins, W., Sima, A., Boucher, O.,
497 Dufresne, J.-L., Nabat, P., and Michou, M.: Effective radiative forcing and adjustments in
498 CMIP6 models, *Atmospheric Chemistry and Physics*, 20, 9591-9618, 2020.

499 Streets, D. G., Fu, J. S., Jang, C. J., Hao, J., He, K., Tang, X., Zhang, Y., Wang, Z., Li, Z., and
500 Zhang, Q.: Air quality during the 2008 Beijing Olympic games, *Atmospheric environment*, 41,
501 480-492, 2007.

502 Su, T., Li, Z., and Kahn, R.: Relationships between the planetary boundary layer height and
503 surface pollutants derived from lidar observations over China: regional pattern and influencing
504 factors, *Atmospheric Chemistry and Physics*, 18, 15921-15935, 2018.

505 Tian, P., Liu, D., Huang, M., Liu, Q., Zhao, D., Ran, L., Deng, Z., Wu, Y., Fu, S., and Bi, K.:
506 The evolution of an aerosol event observed from aircraft in Beijing: An insight into regional
507 pollution transport, *Atmospheric Environment*, 206, 11-20, 2019.

508 Torgo, L.: *Data mining with R: learning with case studies*, Chapman and Hall/CRC2011.

509 Wang, Y., Wang, M., Zhang, R., Ghan, S. J., Lin, Y., Hu, J., Pan, B., Levy, M., Jiang, J. H., and
510 Molina, M. J.: Assessing the effects of anthropogenic aerosols on Pacific storm track using a
511 multiscale global climate model, *Proceedings of the National Academy of Sciences*, 111, 6894-
512 6899, 2014.

513 Wei, J., Li, Z., Lyapustin, A., Sun, L., Peng, Y., Xue, W., Su, T., and Cribb, M.: Reconstructing
514 1-km-resolution high-quality PM_{2.5} data records from 2000 to 2018 in China: spatiotemporal
515 variations and policy implications, *Remote Sensing of Environment*, 252, 112136, 2021.

516 Wei, J., Li, Z., Cribb, M., Huang, W., Xue, W., Sun, L., Guo, J., Peng, Y., Li, J., Lyapustin, A.,
517 Liu, L., Wu, H., and Song, Y.: Improved 1 km resolution PM_{2.5} estimates across China using
518 enhanced space–time extremely randomized trees, *Atmos. Chem. Phys.*, 20, 3273-3289,
519 10.5194/acp-20-3273-2020, 2020.

520 Wu, Y., Liu, D., Wang, X., Li, S., Zhang, J., Qiu, H., Ding, S., Hu, K., Li, W., and Tian, P.:
521 Ambient marine shipping emissions determined by vessel operation mode along the East China
522 Sea, *Science of The Total Environment*, 769, 144713, 2021.

523 Xiao, Q., Chang, H. H., Geng, G., and Liu, Y.: An ensemble machine-learning model to predict
524 historical PM_{2.5} concentrations in China from satellite data, *Environmental science &
525 technology*, 52, 13260-13269, 2018.

526 Xiao, Q., Zheng, Y., Geng, G., Chen, C., Huang, X., Che, H., Zhang, X., He, K., and Zhang,
527 Q.: Separating emission and meteorological contributions to long-term
528 PM_{2.5} trends over eastern China during 2000–2018, *Atmospheric
529 Chemistry and Physics*, 21, 9475-9496, 10.5194/acp-21-9475-2021, 2021.

530 Xu, R., Ye, T., Yue, X., Yang, Z., Yu, W., Zhang, Y., Bell, M. L., Morawska, L., Yu, P., and
531 Zhang, Y.: Global population exposure to landscape fire air pollution from 2000 to 2019, *Nature*,
532 621, 521-529, 2023.

533 Xu, X. and Akhtar, U.: Identification of potential regional sources of atmospheric total gaseous
534 mercury in Windsor, Ontario, Canada using hybrid receptor modeling, *Atmospheric Chemistry
535 and Physics*, 10, 7073-7083, 2010.

536 Xue, T., Zheng, Y., Tong, D., Zheng, B., Li, X., Zhu, T., and Zhang, Q.: Spatiotemporal
537 continuous estimates of PM_{2.5} concentrations in China, 2000–2016: A machine learning
538 method with inputs from satellites, chemical transport model, and ground observations,
539 *Environment international*, 123, 345-357, 2019.

540 Zeng, Z., Gui, K., Wang, Z., Luo, M., Geng, H., Ge, E., An, J., Song, X., Ning, G., and Zhai,
541 S.: Estimating hourly surface PM_{2.5} concentrations across China from high-density
542 meteorological observations by machine learning, *Atmospheric Research*, 254, 105516, 2021.

543 Zhang, J. and Reid, J.: A decadal regional and global trend analysis of the aerosol optical depth
544 using a data-assimilation grade over-water MODIS and Level 2 MISR aerosol products,
545 *Atmospheric Chemistry and Physics*, 10, 10949-10963, 2010.

546 Zhang, L., Wang, T., Lv, M., and Zhang, Q.: On the severe haze in Beijing during January 2013:
547 Unraveling the effects of meteorological anomalies with WRF-Chem, *Atmospheric*
548 *Environment*, 104, 11-21, 2015.

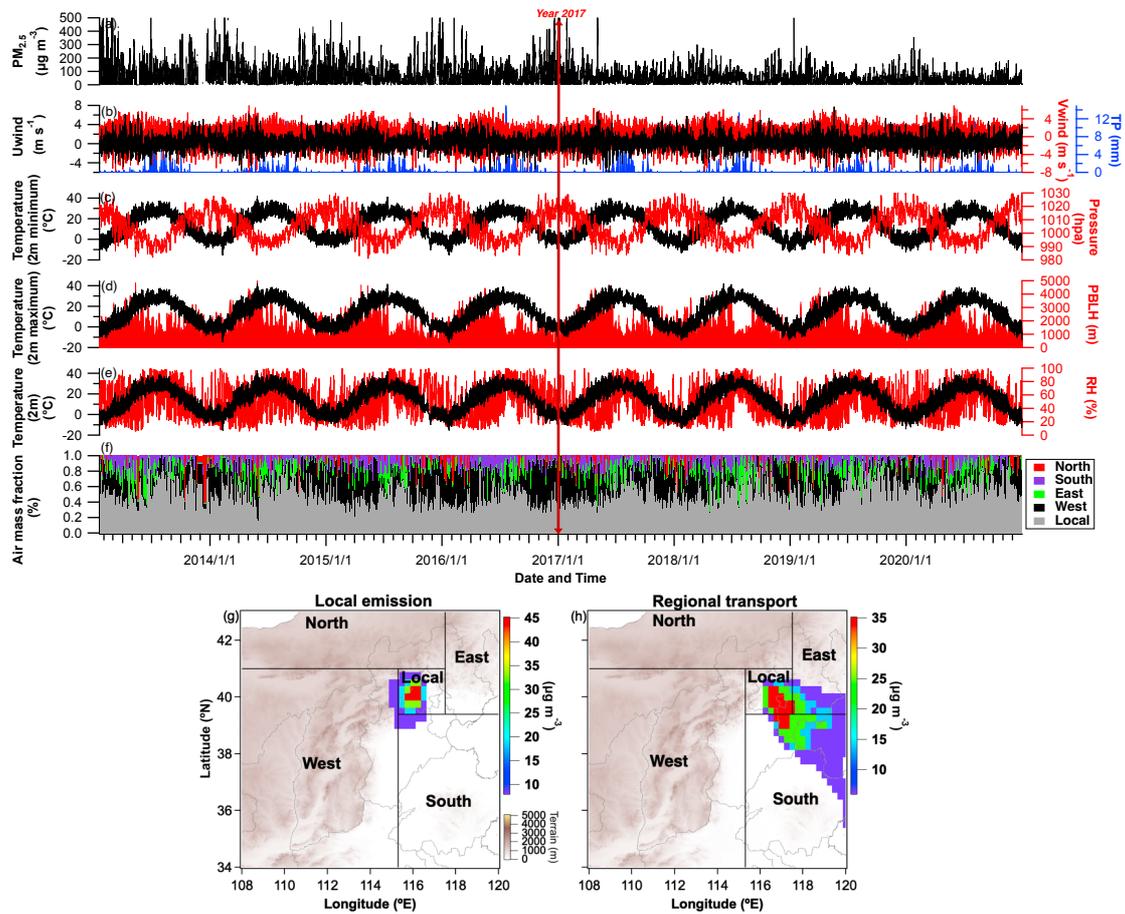
549 Zhang, L., Zhao, T., Gong, S., Kong, S., Tang, L., Liu, D., Wang, Y., Jin, L., Shan, Y., and Tan,
550 C.: Updated emission inventories of power plants in simulating air quality during haze periods
551 over East China, *Atmospheric Chemistry and Physics*, 18, 2065-2079, 2018.

552 Zhang, Q., Wu, S., Wang, X., Sun, B., and Liu, H.: A PM_{2.5} concentration prediction model
553 based on multi-task deep learning for intensive air quality monitoring stations, *Journal of*
554 *Cleaner Production*, 275, 122722, 2020.

555 Zhang, Q., Zheng, Y., Tong, D., Shao, M., Wang, S., Zhang, Y., Xu, X., Wang, J., He, H., and
556 Liu, W.: Drivers of improved PM_{2.5} air quality in China from 2013 to 2017, *Proceedings of*
557 *the National Academy of Sciences*, 116, 24463-24469, 2019.

558

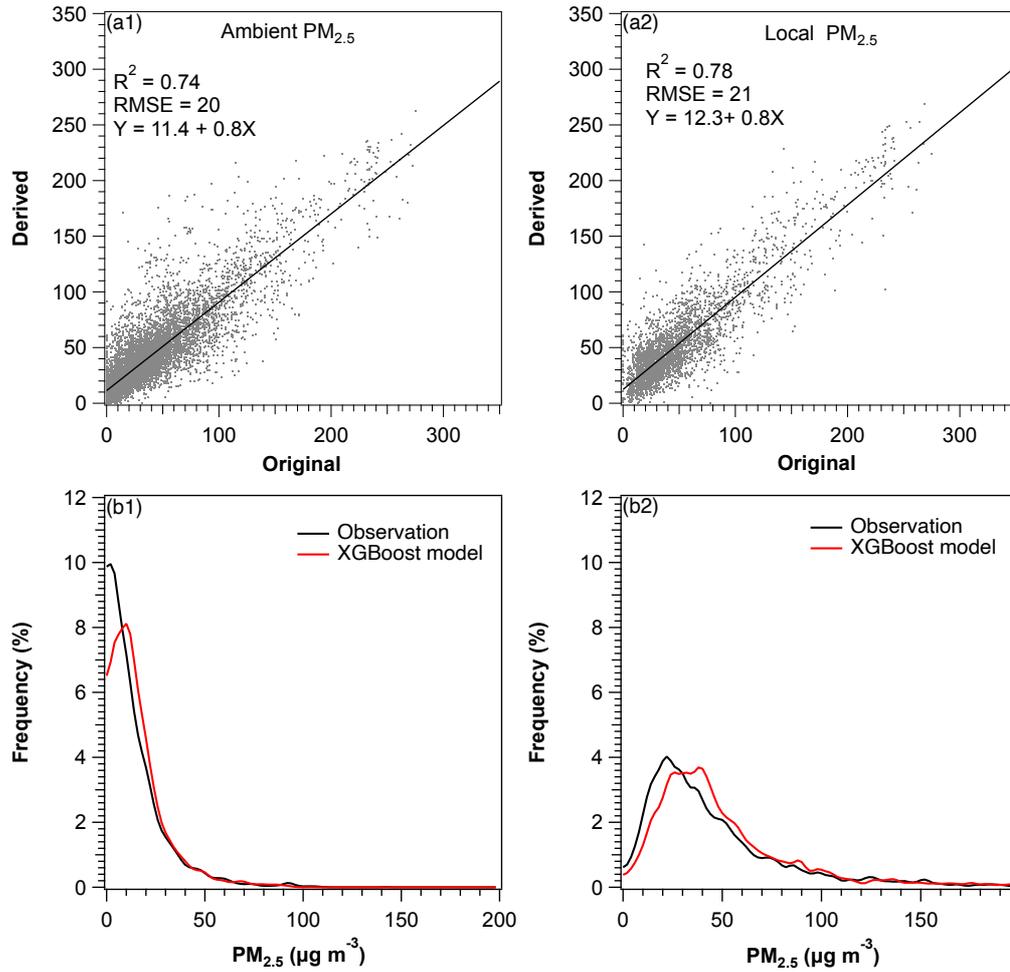
559



561

562 Fig. 1. Temporal evolution of parameters used in the XGBoost model: (a) PM_{2.5}; (b) U-wind,
 563 V-wind, and total precipitation; (c) 2-m minimum temperature and surface pressure; (d) 2-m
 564 maximum temperature and planetary boundary layer height; (e) 2-m temperature and relative
 565 humidity; (f) air mass fraction in contributing sources derived from the Concentration-
 566 Weighted Trajectory (CWT) model for a 1-day backward trajectory. The red vertical line with
 567 arrows indicates the implementation of environmental regulations. Typical examples of the
 568 CWT model analysis are shown for (g) a local emission period (25 August 2013) and (h) a
 569 regional transport period (15 July 2013).

570



571

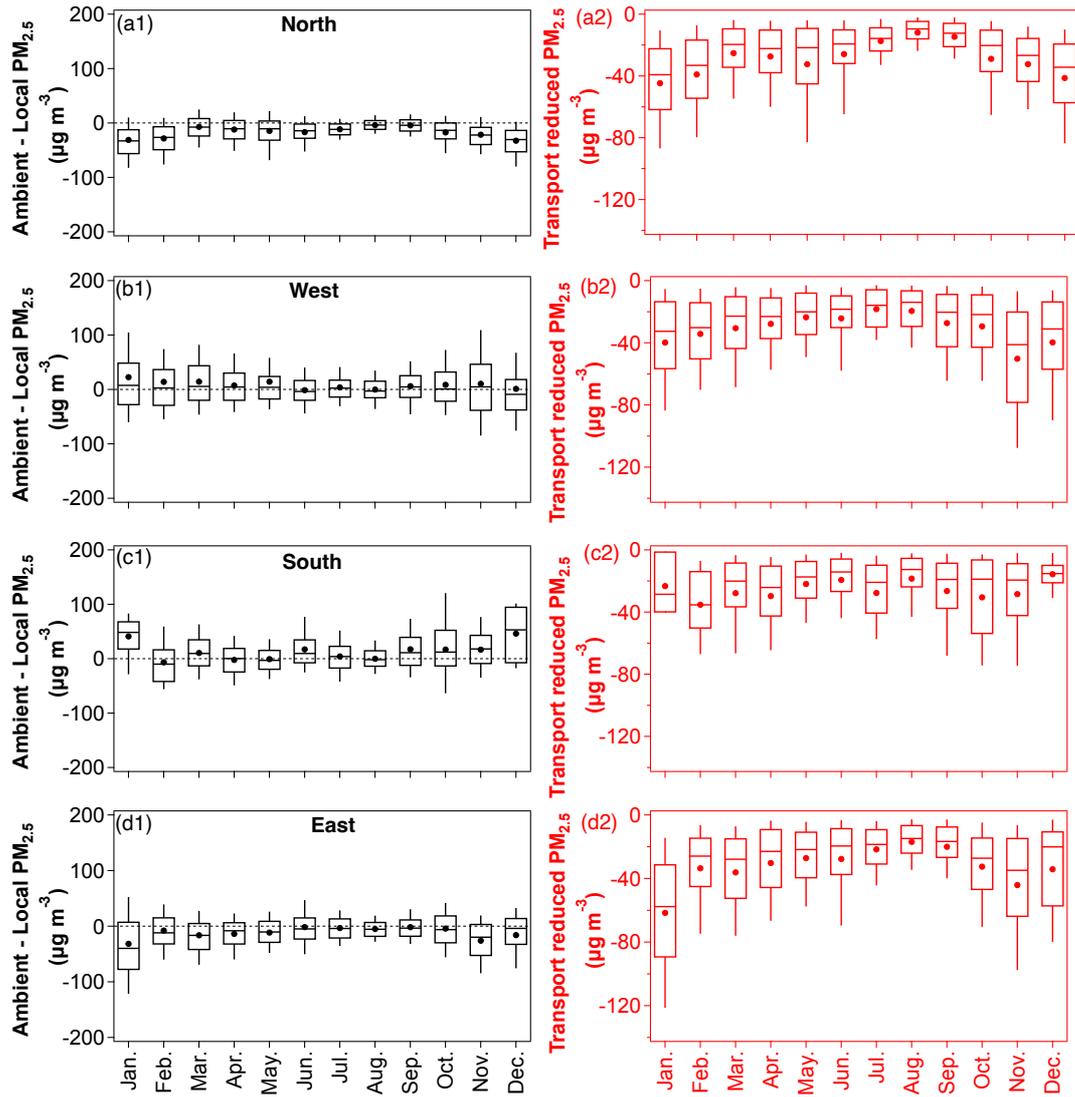
572 Fig. 2. Comparison of XGBoost model estimates and observations for (a1) ambient PM_{2.5} and

573 (a2) local PM_{2.5} using testing samples from 2020. Frequency distributions of PM_{2.5} observations

574 (black lines) and XGBoost model predictions (red lines) for (b1) ambient PM_{2.5} and (b2) local

575 PM_{2.5} using testing samples from 2020.

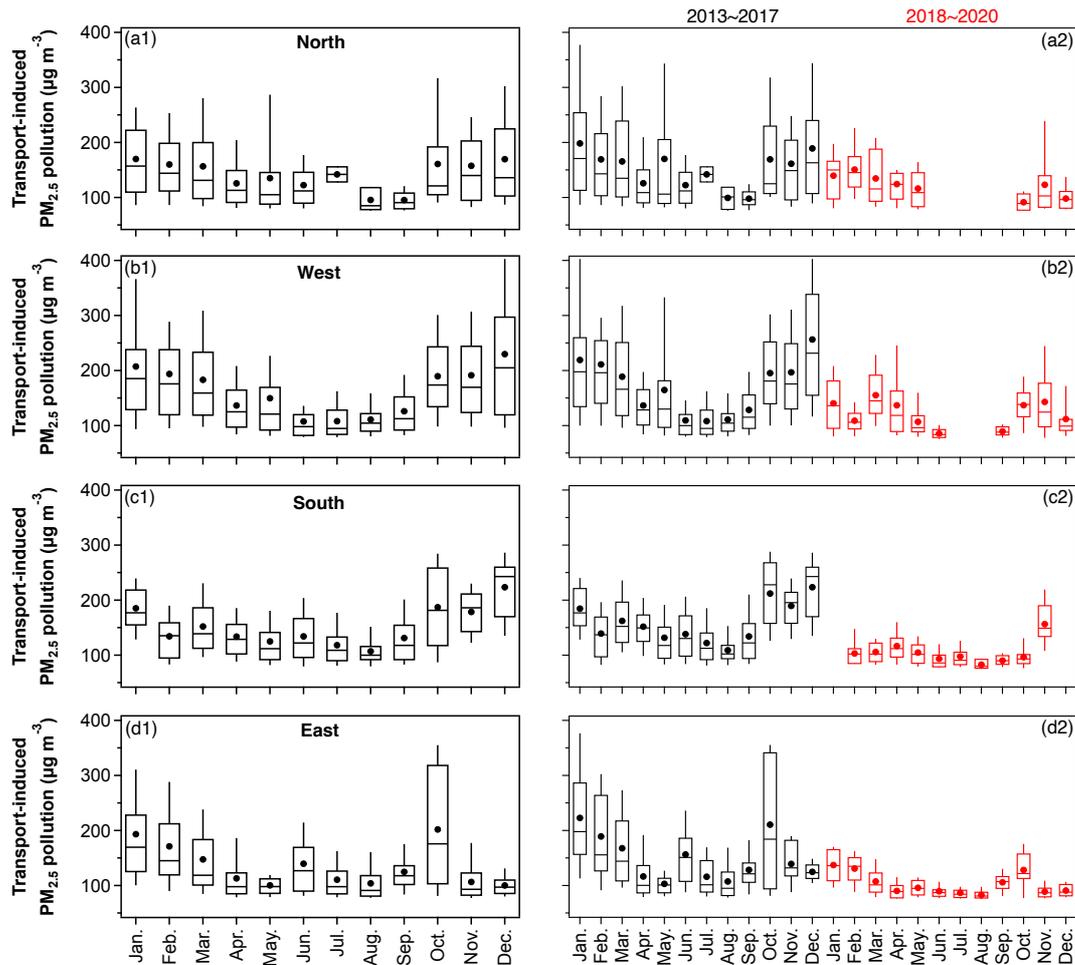
576



577

578 Fig. 3. Monthly variations of the difference between ambient and local $PM_{2.5}$ from the (a1)
 579 North, (b1) West, (c1) South, and (d1) East regions. Right panels show monthly variations of
 580 $PM_{2.5}$ reductions caused by regional transport for the corresponding source regions in the left
 581 panels. The upper and lower boundaries represent the 75th and 25th percentiles, respectively,
 582 while the solid origin represents the average value.

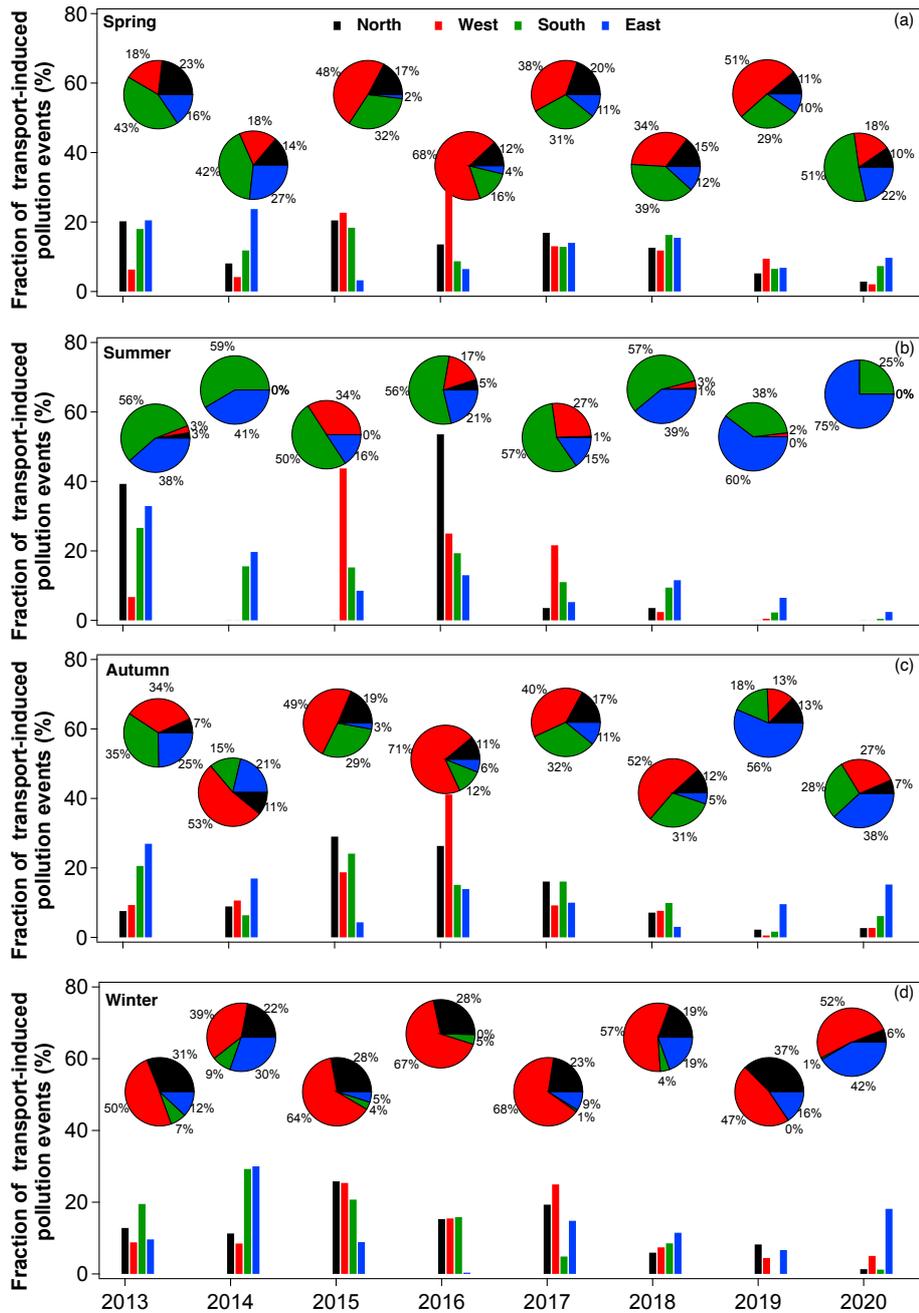
583



584

585 Fig. 4. Monthly variations of transport-induced $PM_{2.5}$ pollution (ambient $PM_{2.5}$ exceeding local
 586 $PM_{2.5}$ and $75 \mu g m^{-3}$) from the (a1) North, (b1) West, (c1) South, and (d1) East regions. Right
 587 panels show monthly variations of transport-induced $PM_{2.5}$ pollution before (black) and after
 588 (red) 2017 for the corresponding source regions in the left panels. The upper and lower
 589 boundaries represent the 75th and 25th percentiles, respectively, while the solid origin represents
 590 the average result.

591



592

593 Fig. 5. Histograms depict the annual fraction of transport-induced pollution events in each
 594 direction relative to the total number of occurrences from 2013 to 2020 during (a) spring, (b)
 595 summer, (c) autumn, and (d) winter. Pie charts illustrate the proportion of transport-induced
 596 pollution events in each direction for each year within the corresponding seasons.

597