

1 **LGHAP: a Long-term Gap-free High-resolution Air Pollutants concentration**  
2 **dataset derived via tensor flow based multimodal data fusion**

3 Kaixu Bai<sup>1,2\*</sup>, Ke Li<sup>1</sup>, Mingliang Ma<sup>3</sup>, Kaitao Li<sup>4</sup>, Zhengqiang Li<sup>4</sup>, Jianping Guo<sup>5\*</sup>,  
4 Ni-Bin Chang<sup>6</sup>, Zhuo Tan<sup>1</sup>, Di Han<sup>1</sup>

5 <sup>1</sup>Key Laboratory of Geographic Information Science (Ministry of Education), School of Geographic Sciences,  
6 East China Normal University, Shanghai 200241, China

7 <sup>2</sup>Institute of Eco-Chongming, 20 Cuiniao Rd., Chongming, Shanghai 202162, China

8 <sup>3</sup>School of Surveying and Geo-Informatics, Shandong Jianzhu University, Jinan 250101, China

9 <sup>4</sup>State Environmental Protection Key Laboratory of Satellite Remote Sensing, Aerospace Information  
10 Research Institute, Chinese Academy of Sciences, Beijing 100101, China

11 <sup>5</sup>State Key Laboratory of Severe Weather, Chinese Academy of Meteorological Sciences, Beijing, China

12 <sup>6</sup>Department of Civil, Environmental, and Construction Engineering, University of Central Florida, Orlando,  
13 FL, USA

14  
15 **\*Correspondence to:** Kaixu Bai ([kxbai@geo.ecnu.edu.cn](mailto:kxbai@geo.ecnu.edu.cn)) and Jianping Guo ([jpguocams@gmail.com](mailto:jpguocams@gmail.com))  
16  
17

18 **Abstract.** Developing a big data analytics framework for generating a Long-term Gap-free High-  
19 resolution Air Pollutants concentration dataset (abbreviated as LGHAP) is of great significance for  
20 environmental management and earth system science analysis. By synergistically integrating  
21 multimodal aerosol data acquired from diverse sources via a tensor flow based data fusion method, a  
22 gap-free aerosol optical depth (AOD) dataset with daily 1-km resolution covering the period of 2000–  
23 2020 in China was generated. Specifically, data gaps in daily AOD imageries from MODIS aboard  
24 Terra were reconstructed based on a set of AOD data tensors acquired from [diverse](#) satellites,  
25 numerical analysis, and *in situ* air quality [measurements](#) via integrative efforts of spatial pattern  
26 recognition for high dimensional gridded image analysis and knowledge transfer in statistical data  
27 mining. To our knowledge, this is the first long-term gap-free high resolution AOD dataset in China,  
28 from which spatially contiguous PM<sub>2.5</sub> and PM<sub>10</sub> concentrations were [then](#) estimated using an  
29 ensemble learning approach. Ground validation results indicate that the LGHAP AOD data are in a  
30 good agreement with *in situ* AOD observations from AERONET, with R of 0.91 and RMSE equaling  
31 to 0.21. Meanwhile, PM<sub>2.5</sub> and PM<sub>10</sub> estimations also agreed well with ground measurements, with R  
32 of 0.95 and 0.94 and RMSE of 12.03 and 19.56 μg m<sup>-3</sup>, respectively. [The](#) LGHAP provides a suite of  
33 long-term gap free gridded maps with high-resolution to better examine aerosol changes in China over  
34 the past two decades, from which three [major](#) variation periods of haze pollution were revealed in  
35 China. Additionally, the proportion of population exposed to unhealthy PM<sub>2.5</sub> was increased from  
36 50.60% in 2000 to 63.81% in 2014 across China, which was then [reduced drastically](#) to 34.03% in  
37 2020. Overall, the generated LGHAP [dataset](#) has a great potential to trigger multidisciplinary  
38 applications in earth observations, climate change, public health, ecosystem assessment, and  
39 environmental management. The daily resolution AOD, PM<sub>2.5</sub>, and PM<sub>10</sub> datasets [are](#) publicly  
40 [available](#) at <https://doi.org/10.5281/zenodo.5652257> (Bai et al., 2021a),  
41 <https://doi.org/10.5281/zenodo.5652265> (Bai et al., 2021b), and  
42 <https://doi.org/10.5281/zenodo.5652263> (Bai et al., 2021c), respectively. [Monthly](#) and annual datasets  
43 can be [acquired from](#) <https://doi.org/10.5281/zenodo.5655797> (Bai et al., 2021d) and  
44 <https://doi.org/10.5281/zenodo.5655807> (Bai et al., 2021e), respectively. Python, Matlab, R, and IDL  
45 codes were also provided to help users read and visualize these data.

46 **Keywords:** Aerosol optical depth; Particulate matter; Gap filling; Big data analytics; Multimodal data  
47 fusion

Deleted: ¶

Deleted: data

Deleted: Overall, t

Deleted: distinct

Deleted: drastically

Deleted: aerosol

Deleted: can be

Deleted: accessed

Deleted: Meanwhile,

Deleted: m

Deleted: mean

Deleted: found

Deleted: at

61 **1 Introduction**

62 Atmospheric aerosols not only impact regional climate by changing the Earth radiation budget  
63 but significantly influence air quality at the ground level (Fuzzi et al., 2015; Gao et al., 2018; Shen et  
64 al., 2020; Sun et al., 2015; [Yang et al., 2020](#); [Zheng et al., 2020](#)). Monitoring aerosol loading in the  
65 atmosphere is thus of great significance for climate change attribution and haze pollution assessment.  
66 Aerosol optical depth (AOD), an [indicator](#) of [aerosol bulks](#) distributed within a column of air from the  
67 Earth's surface to the top of the atmosphere, has been monitored for decades to [map global](#) aerosol  
68 loading in the atmosphere. Compared with sparsely [and unevenly](#) distributed ground-based aerosol  
69 monitoring stations (e.g., AERONET), satellite instruments can [map AOD with vaster](#) spatial coverage  
70 [at even sub-hourly](#) sampling frequency (e.g., [geostationary satellite](#)). An overview of sensors,  
71 algorithms, and AOD datasets that are widely used [in the community](#) can be found in the literature  
72 such as Sogacheva et al. (2020) and Wei et al. (2020).

73 Due to negative impacts of bright surface (e.g., snow cover) and clouds, as well as algorithmic  
74 restrictions, satellite AOD retrievals often suffer from extensive data gaps, significantly reducing the  
75 downstream application potential such as mapping particulate matter (PM) concentrations at the  
76 ground surface (e.g., Bai et al., 2019a; Wei et al., 2021a). Also, [excessive](#) data gaps in AOD imageries  
77 may result in large uncertainty when assessing aerosol impacts on weather and climate (Guo et al.,  
78 2017; Li et al., 2019; Zhao et al., 2020; [Zheng et al., 2018](#)). Over the years, [versatile](#) gap filling methods  
79 have been developed (e.g., Bai et al., 2016, 2020b; Chang et al., 2015). Nonetheless, filling data gaps  
80 in satellite-based AOD [retrievals](#) is still [challenging](#) due to extraordinary nonrandom missing values  
81 and high aerosol dynamics in space and time.

82 [Wei et al. \(2020\)](#) provided a short review of methods that have been frequently applied to deal  
83 with data gaps in AOD products. In general, merging AOD data acquired from diverse instruments  
84 and/or platforms is the most popular approach to improve AOD spatial coverage (Sogacheva et al.,  
85 2020). Statistical methods such as linear regression (Bai et al., 2019a; Wang et al., 2019; Zhang et al.,  
86 2017), inversed variance weighting (Chen et al., 2018; Ma et al., 2016; Sogacheva et al., 2020), and  
87 maximum likelihood estimate (Xu et al., 2015), are often applied to account for systematic bias among  
88 different datasets. Data fusion methods such as Bayesian maximum entropy could be applied to blend  
89 AOD products with different resolutions (Tang et al., 2016; Wei et al., 2021b). Another way is to  
90 reconstruct missing AOD values using either neighboring observations in space and time or external

Field Code Changed

Deleted: measure

Deleted: aerosols

Deleted: quantify

Deleted: provide better

Deleted: observations because of

Deleted: and high

Field Code Changed

Deleted: many

Deleted: products

Deleted: a challenge

Field Code Changed

Deleted: a

101 data sources such as AOD simulations from numerical models (Li et al., 2020; Xiao et al., 2017), even  
102 meteorological factors (Bi et al., 2018).

Deleted: a..., even simply ... [1]  
Field Code Changed

103 Although there exist a variety of gap filling methods, spatially gap free AOD datasets are still  
104 rare, particularly high-resolution AOD datasets from satellites, significantly limiting downstream  
105 applications such as PM<sub>x</sub> concentration mapping. In spite of versatile PM<sub>2.5</sub> concentration prediction  
106 models (e.g., Di et al., 2019; Fang et al., 2016; Hu et al., 2014; Li et al., 2016; Lin et al., 2016; Liu et  
107 al., 2009; Wang et al., 2021a), to date, there are few publicly accessible PM<sub>x</sub> concentration datasets  
108 that can be used to examine haze pollution variations regionally and globally. Several typical datasets,  
109 e.g., the one generated by the Dalhousie University (van Donkelaar et al., 2010, 2016), CHAP (Wei et  
110 al., 2021a), and TAP (Geng et al., 2021), have been widely applied to advance our understanding on  
111 aerosol impacts across China and globe. However, these datasets more or less still suffer from  
112 drawbacks in spatial and/or temporal resolution, spatial coverage, and data accuracy. To meet the  
113 contemporary needs, Zhang et al. (2021) provided a more comprehensive review of the widely used  
114 PM<sub>x</sub> concentration mapping approaches. With a thorough review for PM<sub>2.5</sub> concentration mapping  
115 techniques, an optimal full-coverage PM<sub>2.5</sub> modeling scheme was proposed, in which diverse aerosol  
116 datasets were fused toward a full-coverage AOD map based on a multi-modal approach (Bai et al.,  
117 2022). In parallel with these efforts, some attempted to improve AOD data coverage over space with  
118 high accuracy by merging AODs observed at adjacent times directly (Li et al., 2022).

Deleted: many versatile... variety of gap filling methods, spatially gap free AOD datasets are always ...till rare, particularly satellite-based ...igh-resolution AOD datasets from satellites, resulting in s... significantly limiting in ... [2]  
Field Code Changed

Deleted: He et al., 2020; ...u et al., 2014; Li et al., 2018b...016, 2016... Lin et al., 2016; Liu et al., 2009; Ma et al., 2014; ... [3]  
Formatted: Subscript

Deleted: Aa ...everal typical...ypical ... [4]  
Field Code Changed

Deleted: 2019a  
Field Code Changed

Deleted: a..., have been demonstrated the global effort to...idely applied to elevate ...dvance our understanding on earth system science research ... [5]  
Formatted: Subscript

Deleted: ¶ ... [6]

119 With such prior knowledge, the current study developed a big data analytics framework for  
120 generating a Long-term Gap-free High-resolution Air Pollutants concentration dataset (abbreviated as  
121 LGHAP hereafter), aiming at providing gap-free AOD, PM<sub>2.5</sub> and PM<sub>10</sub> concentration data with a daily  
122 1-km resolution in China for the period of 2000 to 2020. Toward such a goal, multimodal aerosol data  
123 acquired from diverse sources including satellites, ground stations and numerical models were  
124 synergistically integrated via the higher order singular value decomposition (HOSVD) to form a tensor  
125 flow based data fusion framework in the current study. Full coverage PM<sub>2.5</sub> and PM<sub>10</sub> concentration  
126 data were then estimated on the basis of the gap-filled AOD dataset. This 21-year-long gap-free high  
127 resolution (daily/1km) aerosol dataset was then compared against ground-based AOD and PM<sub>x</sub>  
128 observations to validate the data accuracy of each product, particularly their performance in spatial  
129 pattern recognition and temporal trend assessment. These advances endorsed a better assessment of

Deleted: Given ...ith such prior knowledge, the current study developed a big data analytics framework for generating a Long-term Gap-free High-resolution Air Pollutants concentration dataset (abbreviated as LGHAP hereafter), aiming at providing gap-free AOD, PM<sub>2.5</sub> and PM<sub>10</sub> concentration data with a daily 1-km resolution in China from ...or the period of 2000 to 2020. Toward such a goal, mM...ltimodal aerosol data acquired from diverse sources including satellites, ground stations and numerical models were synergistically integrated via the higher order singular value decomposition (HOSVD) to form a tensor flow based data fusion method ... [7]  
Formatted: Subscript

Deleted: evaluate ...alidate the data accuracy of each product, particularly their performance in spatial pattern recognition and temporal trend assessment. ...hese advances endorsed a led to...etter explor...ssessment ofe ... [8]

194 long-term variability of haze pollution in China as well as the corresponding population exposure over  
 195 the past two decades.

196 **2 Data sources**

197 Table 1 provides a brief summary of the multisource datasets used in this study to generate the  
 198 LGHAP dataset. As shown, six satellite-based AOD products, five numerical simulations of AOD and  
 199 aerosol components, eleven meteorological factors, six datasets of ground-based AOD and air  
 200 pollutants concentration measurements, as well as a set of land cover, topographic and socioeconomic  
 201 parameters, were employed. Descriptions of these datasets are given in the following subsections.

202 **Table 1.** Summary of the data sources used in this study to generate gap free high resolution AOD  
 203 and PM<sub>x</sub> concentration datasets.

Category	Source product	Time range	Temporal resolution	Spatial resolution
AOD	Terra/MODIS	2000–2020	daily	1 km
	Aqua/MODIS	2002–2020	daily	1 km
	Terra/MISR	2000–2020	daily	4.4 km
	Suomi-NPP/VIIRS	2012–2020	daily	5 km
	Envisat/AATSR	2000–2012	daily	10 km
	PARASOL/POLDER	2005–2013	daily	10 km
	MERRA-2	2000–2020	hourly	0.5°×0.625°
Meteorology	AERONET	2000–2020	hourly	point
	Air temperature		hourly	0.25°
	U/V component of wind		hourly	0.25°
	Relative humidity		hourly	0.25°
	Surface pressure		hourly	0.25°
	Boundary layer height	2000–2020	hourly	0.25°
	Total column water vapor		hourly	0.25°
	Surface solar radiation downwards		hourly	0.25°
	Instantaneous moisture flux		hourly	0.25°
	Visibility	2000–2013	3-hour	point
Air quality	PM <sub>2.5</sub> , PM <sub>10</sub> , SO <sub>2</sub> , NO <sub>2</sub>	2014–2020	hourly	point
Population	WorldPop	2000–2020	annual	1 km
Elevation	DEM	2000	/	30 m
Land Cover	CLCD	2000–2019	annual	30 m
	GLOBELAND	2020	annual	30 m
NDVI	Terra/MODIS	2000–2020	monthly	1 km
Aerosol component	MERRA-2	2000–2020	hourly	0.5°×0.625°

Deleted: the

Deleted: and

Deleted: to haze pollution in China

Deleted: by taking advantage of the LGHAP dataset

Deleted: summarizes

Deleted: help

Deleted: product

Deleted: quality datasets

Deleted: five

Formatted: Subscript

213 **2.1 Gridded aerosol products**

214 In many previous studies, coarse AOD and/or aerosol components simulations acquired from  
215 numerical models were oftentimes used as the primary data source to help derive full-coverage AOD  
216 and/or PM<sub>2.5</sub> concentration maps (e.g., Park et al., 2020; Wang et al., 2021b). However, due to the lack  
217 of high accuracy near real-time emission inventory, simulated AOD and/or aerosol components are  
218 often prone to large uncertainty, which could be inevitably introduced to the final PM<sub>2.5</sub> estimations if  
219 no observational data are applied for possible bias correction. In such a research context, here we used  
220 six satellite-based AOD products with a relatively long temporal coverage (>5 years) to help better  
221 reconstruct historical AOD variations over space and time, though geostationary satellites can provide  
222 AOD observations at even hourly resolution. The reasons are twofold. On the one hand, the operational  
223 AOD product from the recent Chinese FY-4 satellite is still unavailable. On the other hand, AOD  
224 product from Hamawari-8 cannot provide observations in the northwest region of China.

225 The latest AOD product derived from the MODerate-resolution Imaging Spectroradiometer  
226 (MODIS) onboard Terra using the multiangle implementation of atmospheric correction (MAIAC)  
227 algorithm (Lyapustin et al., 2011, 2018), was hereby used as the baseline dataset for the generation of  
228 gap free AOD maps. This AOD product has not only a finer spatial resolution (1 km) but a comparable  
229 and even better accuracy, when comparing with those derived from the Dark Target and Deep Blue  
230 algorithms (Goldberg et al., 2019; Lyapustin et al., 2018). In addition, AOD products derived from  
231 MODIS onboard Aqua, the Multi-angle Imaging SpectroRadiometer (MISR) onboard Terra, Visible  
232 Infrared Imaging Radiometer Suite (VIIRS) onboard Suomi-NPP, Advanced Along-Track Scanning  
233 Radiometer (AATSR) onboard Envisat and POLarization and Directionality of the Earth's  
234 Reflectances (POLDER) onboard PARASOL, were also employed. The ultimate goal was to reduce  
235 the bias level in the final full-coverage AOD product by providing observational AODs as much as  
236 possible. Accuracies of these AOD products have been extensively validated in previous studies, e.g.,  
237 de Leeuw et al. (2018), Xiao et al. (2016), Wei et al. (2019b), Che et al. (2019), to name a few. A brief  
238 description of these satellite-based AOD products can be found in Text S1 in the supplementary  
239 information.

240 In addition to satellite-based AOD products, numerically simulated aerosol diagnostics from  
241 MERRA-2, including AOD and aerosol components such as black carbon, organic carbon, dust and  
242 sulfate, were also applied to help reconstruct missing AOD information and to predict PM<sub>2.5</sub> and PM<sub>10</sub>

Deleted: -term

Deleted: .

Deleted: full-coverage

Field Code Changed

Deleted: ; Xiao et al., 2017b

Deleted: utilized

Deleted: imagery

Deleted: better spatial coverage of

Deleted:

Deleted: , in the current study

252 concentrations at the ground level. The aerosol components were used here as a proxy of emission  
253 inventory when predicting  $PM_x$  concentrations. Big data analytics procedures applied to these datasets  
254 will be described in section 3.

Formatted: Subscript

## 255 2.2 *In situ* AOD and air quality measurements

256 AOD observations from Aerosol Robotic Network (AERONET) were ~~hereby used as the ground~~  
257 ~~truth~~ to evaluate the ~~data~~ accuracy of the generated ~~gap free~~ AOD product, as well as the learning  
258 target to infer AOD from air pollutants concentration and atmospheric visibility. Considering few valid  
259 data were provided in the Level 2.0 dataset, here we used the Level 1.5 AOD data to guarantee adequate  
260 *in situ* AOD data coverage in space and time. To validate the gridded AOD products in this study, each  
261 *in situ* AOD observation was registered with the gridded mean AOD over a 50×50 km window.

Deleted: prediction

Deleted: full-coverage (

Deleted: )

262 Near-surface air pollutants concentrations including  $PM_{2.5}$ ,  $PM_{10}$ ,  $NO_2$ , and  $SO_2$  that were  
263 sampled at state-controlled monitoring sites were also applied, not only to help establish machine-  
264 learned regression models for  $PM_x$  prediction ( $PM_{2.5}$  and  $PM_{10}$ ), but to infer AOD over air quality  
265 monitoring sites given their dense distributions across China. The gauged air pollutants concentration  
266 data have been released online on an hourly basis by the China National Environment Monitoring  
267 Center since the late 2013. For quality control, outliers were first detected and removed from each  
268 pollutant dataset by following the criteria used in our previous study (Bai et al., 2020a). The missing  
269 values were then reconstructed using the diurnal cycle constrained empirical orthogonal function  
270 (DCCEOF) method proposed in Bai et al. (2020b).

Formatted: Subscript

271 The 3-hour resolution atmospheric visibility data acquired from 4,052 weather stations were  
272 employed to help generate gap free AOD maps before 2014, at which *in situ* air quality measurements  
273 were not available. Previous studies have attempted to predict  $PM_{2.5}$  concentration from atmospheric  
274 visibility data with good accuracies (Liu et al., 2017), indicative of a great potential for estimating  
275 AOD. Specifically, visibility data were used as an important predictor for site-specific AOD prediction,  
276 and the resulting AOD predictions were then used as a critical prior information for reconstructing  
277 AOD distributions over space, especially over those regions without satellite AOD observations. ~~Given~~  
278 ~~the availability of abundant air quality measurements and the fact that~~ automatic visibility sensors have  
279 been widely used ~~across China~~ since 2014, ~~atmospheric visibility data after 2014~~ were ~~thereby~~  
280 excluded to guarantee the data consistency (Li et al., 2018a). For quality control, the consistency of

Deleted: Since

Deleted: at many sites

Deleted: those

287 visibility data was examined using an outlier detection method, i.e., the annual mean should not exceed  
288 3 times the standard deviation of data over a 5-year time window (Zhang et al., 2020). Those with  
289 apparent jumps and drifts in visibility time series were excluded. Meanwhile, visibility data on  
290 rainstorm and foggy days were eliminated as well.

### 291 2.3 Auxiliary data

292 As shown in Table 1, eleven meteorological factors, including air temperature at the near surface,  
293 wind speed and direction, relative humidity, surface pressure, boundary layer height, total column  
294 water vapor, downwards solar radiation, and instantaneous moisture flux, were used to help **resolve**  
295 **nonlinear relationships between PM<sub>2.5</sub> and AOD**, as well as to downscale AOD from MERRA-2. These  
296 data were acquired from the fifth generation ECMWF atmospheric reanalysis (ERA-5), and the first  
297 three factors were extracted at the levels of not only ground surface but 850 hpa and 500 hpa so as to  
298 indicate the vertical structure of the atmosphere. Additionally, population data from WorldPop, land  
299 cover from CLCD during 2000 to 2019 (Yang and Huang, 2021) and GLOBELAND 30 in 2020 (Chen  
300 et al., 2014), elevation data from the Global Digital Elevation Model (GDEM) version 2, as well as  
301 monthly composited 1-km normalized difference vegetation index (NDVI) from MODIS, were  
302 employed to **resolve** the socioeconomic and ecological contributions to haze pollutions. Properties of  
303 these datasets can be found in Table 1, and datasets with a finer resolution were upscaled to 0.01° via  
304 a cubic interpolation method.

### 305 3 Methodology

306 **Toward the generation of LGHAP aerosol datasets** to advance environment management and  
307 earth system science analysis, **here we** developed a big data analytics framework **via** a seamless  
308 integration of the tensor flow based multimodal data fusion with ensemble learning based PM<sub>2.5</sub>  
309 concentration estimation. **The** proposed method transformed a set of data tensors of AOD and other  
310 related datasets such as air pollutants concentration and atmospheric visibility that were acquired from  
311 diversified sensors or platforms via integrative efforts of spatial pattern recognition for high  
312 dimensional gridded data analysis toward data fusion and multiresolution image analysis, as well as  
313 knowledge transfer in statistical data mining. The proposed method consists of three major procedures  
314 in general, including multisensory data homogenization, tensor flow based AOD reconstruction, and

Deleted: infer

Deleted: 2.5

Deleted: and PM<sub>10</sub> from

Deleted: indicate

Deleted: T

Deleted: the current study

Deleted: for generating long-term gap free aerosol spatiotemporal datasets and demonstrated its applications in China. Such big data analytics was constructed

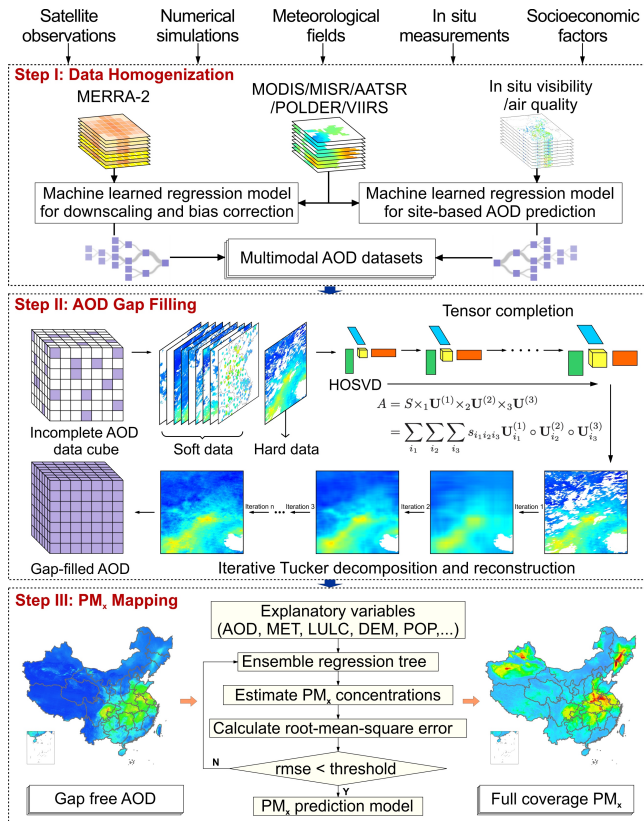
Formatted: Subscript

Deleted: When generating this dataset, t



325 ensemble learning for  $PM_{2.5}$  concentration estimation. The analytical framework of the big data  
 326 analytics is depicted in Figure 1 and described in details in the following subsections.

Formatted: Subscript



327  
 328 **Figure 1.** Flowchart of the proposed big data analytics framework for generating a long-term gap-free  
 329 high-resolution air pollutants concentration dataset (LGHAP), taking aerosol optical depth (AOD) and  
 330  $PM_{2.5}$  ( $PM_{2.5}$  and  $PM_{10}$ ) concentration in China as illustration. HOSVD is an acronym of high order  
 331 singular value decomposition. MET, LULC, DEM, and POP denote variables of meteorology, land  
 332 use/land cover, digit elevation model, and population, respectively.

Formatted: Subscript

Formatted: Subscript

Formatted: Subscript

333 **3.1 Multisensory data homogenization**

334 Since a set of aerosol products with different types, resolution, and accuracies were applied to  
 335 support the reconstruction of gap-free AOD imageries, harmonizing cross-platform biases and scale

Deleted: generation

Deleted: mission

338 differences between these diversified datasets is crucial to multisensory data integration. In this study,  
339 machine-learned regression models were established to harmonize these heterogeneous aerosol  
340 datasets. A baseline dataset was first selected to be used as the learning target while other datasets  
341 were calibrated to the level of baseline dataset to make them comparable. Given finer resolution and  
342 higher proportion of data coverage in space and time, the MAIAC AOD product from Terra (AOD<sub>Terra</sub>)  
343 was selected as the baseline dataset. Consequently, six machine-learned regression models were  
344 established between AOD<sub>Terra</sub> and each gridded AOD product (i.e., five satellite-based AOD products  
345 plus MERRA-2 AOD simulations) using the random forest method. Meteorological factors (MET),  
346 land cover types (LULC), topographic (DEM) and population (POP) were used as covariates to help  
347 downscale these multimodal AOD products to have a resolution same as AOD<sub>Terra</sub> while accounting  
348 for cross-mission biases arising from temporal and algorithmic differences.

Deleted: thus of critical importance

Deleted: facilitate

349 Considering data gaps are extensive in satellite AOD products, especially over regions with vast  
350 cloud cover, providing prior AOD information over such region is thus of great value in support of the  
351 reconstruction of missing AOD values. As indicated in our recent studies, AOD can be accurately  
352 predicted from ground measured air pollutants concentration, showing an accuracy even over some  
353 satellite AOD retrievals (Li et al., 2021; Bai et al., 2021). To support AOD reconstruction over regions  
354 with less or even without valid satellite AOD observations, we attempted to infer AOD over air quality  
355 monitoring sites from in situ air pollutants concentration measurements via a machine learning  
356 approach. Similarly, machine-learned regression models were established using random forest by  
357 taking AOD<sub>Terra</sub> as the learning target while ground measured air pollutants concentration,  
358 meteorological factors, land cover, and terrain information, were used conjunctively as predictors.

Deleted: -based

Deleted: thick

Deleted: without

Deleted: an ensemble

Formatted: Font: Italic

359 The transformation of ground measured air pollutants concentration data to AOD allows for  
360 providing external observational AOD data to supplement satellite observations, especially over  
361 regions suffering from significant data gaps. Since air pollutants concentration data were not available  
362 before 2013, atmospheric visibility data sampled at dense weather stations were hereby used as an  
363 alternative for site-based AOD prediction, by applying a similar prediction model as used above for  
364 air pollutants concentration. Figure S1 show the ground-based validation results of AOD inferred from  
365 atmospheric visibility and air pollutants concentration, indicative of a generally good accuracy of these  
366 inferred AOD values. All efforts led to aggregate a set of multimodal aerosol data with different  
367 properties for multisensory data fusion toward gap free AOD mapping as the next step.

Deleted: empowers us to

Deleted: e

Deleted: generating full-coverage (

Deleted: )

379 **3.2 Tensor flow based AOD reconstruction**

380 The core of generating full coverage AOD imageries is to fill in data gaps in AOD<sub>Terra</sub>. Previous  
381 studies have demonstrated that merging satellite AOD retrievals at adjacent time steps can help  
382 improve the observational AOD coverage at each single snapshot, while the involvement of numerical  
383 AOD simulations can help bridge AOD data gaps (Li et al., 2022; Bai et al., 2022). In this study, a  
384 tensor completion method was particularly designed and applied to fulfil the gap filling in AOD<sub>Terra</sub>.  
385 Specifically, the incomplete AOD<sub>Terra</sub> imageries were deemed as the hard data (true AOD state) while  
386 other AOD datasets (e.g., the downscaled AOD datasets and site-specific AOD predictions inferred  
387 from air pollutants concentration and atmospheric visibility) were used as the soft data (complementary  
388 data) to help reconstruct AOD distribution in AOD<sub>Terra</sub> via tensor flow based pattern recognition.  
389 Detailed procedures for gap filling are outlined as follows.

390 3.2.1 Initial AOD tensor construction

391 Due to extensive data gaps in satellite-based AOD retrievals, it is insufficient to reconstruct all  
392 missing AOD information in AOD<sub>Terra</sub> for a given date by simply merging the harmonized satellite-  
393 based AOD data synchronously. To fulfill AOD gap filling, the tensor completion method was thus  
394 applied to synergistically integrate AOD acquired from diverse sources. Consequently, creating the  
395 data tensor of AOD is of critical importance. In this study, the data tensor of AOD was constructed by  
396 incorporating not only observational AOD from both satellites and those inferred from *in situ* air  
397 quality indicators on the same date, but also historical AOD retrievals from MODIS instruments  
398 (AOD<sub>Terra</sub> and AOD<sub>Aqua</sub>) and part of data from the downscaled MERRA-2 AOD (denoted as AOD<sub>M2</sub>  
399 hereafter). The latter two were applied to provide knowledge of AOD distributions over space to guide  
400 the reconstruction of missing values in AOD<sub>Terra</sub>.

401 For the screening of historical observations resembling AOD<sub>Terra</sub> distribution on the given date  
402 to be reconstructed, AOD<sub>M2</sub> was used in concert with AOD<sub>Terra</sub> and site-based AOD estimations to  
403 identify similar imageries. Toward this goal, site-specific AOD estimations and 5% randomly selected  
404 downscaled AOD<sub>M2</sub> data were merged directly with valid AOD<sub>Terra</sub> to form a new image on each date.  
405 Subsequently, correlations and biases were estimated between AOD<sub>Terra</sub> on the given date to be  
406 reconstructed and each newly merged historical AOD<sub>Terra</sub> image. To avoid the inclusion of imageries  
407 with distinct variation patterns, only those closely resembling AOD<sub>Terra</sub> on the date to be reconstructed

Deleted: derived

Deleted: fill in

Deleted: 2021

Deleted: 2021

Deleted: s

Deleted: using

Deleted: newly developed

Deleted: diversified

Deleted: observations

Formatted: Subscript

Formatted: Subscript

Deleted: select

Deleted: Specifically,

Deleted: ,

Deleted: which was then used to find similar

Deleted: maps

Deleted:

Deleted: AOD distribution in the composite image

Deleted: given

425 were **finally retained** in terms of their correlations and biases subject to a threshold of  $R > 0.7$  and  
 426  $RMSE < 0.2$ . Once sufficient historical imageries were obtained, the data tensor of AOD was  
 427 constructed by compiling the observed AOD imageries on the given date with historical imageries to  
 428 a three-dimension data array  $\mathbf{A} \in \mathbf{R}^{N_1 \times N_2 \times N_3}$  (composed of  $N_3$  images with a size of  $N_1 \times N_2$ ).  
 429 **Considering satellite AOD retrievals suffer from extensive data gaps, we injected data values of site-**  
 430 **specific AOD estimations and 1% randomly selected downscaled AOD<sub>M2</sub> data directly onto grids**  
 431 **where AOD<sub>Terra</sub> values missed on each specific date as prior knowledge. This not only accelerates**  
 432 **convergence speed during the reconstruction process but avoids large reconstruction errors over**  
 433 **regions with tremendous data gaps in satellite observed AOD imageries.**

### 434 3.2.2 Gap filling via tensor completion

435 Previous studies have well demonstrated the good performance of matrix decomposition  
 436 methods such as empirical orthogonal function and singular value decomposition (SVD) for missing  
 437 value imputation (Bai et al., 2020b; Beckers and Rixen, 2003; Folch-Fortuny et al., 2015). However,  
 438 these methods can only work on two-dimension matrix mathematically, **namely** the matrix domain. To  
 439 integrate spatial features of AOD revealed by datasets to generate a smooth AOD distribution with  
 440 complete coverage, in this study, the HOSVD, a specific orthogonal Tucker decomposition, was  
 441 applied. More detailed descriptions to HOSVD can be found in the literature such as Sun et al. (2021),  
 442 Tucker (1966), Kolda and Bader (2009), **and** Sidiropoulos et al. (2017).

443 In Table 2, we provided a stepwise description of the algorithm used to fill data gaps in AOD<sub>Terra</sub>  
 444 by integrating AOD features recognized in different imageries as the data tensor of AOD via HOSVD.  
 445 To initiate the tensor decomposition, grids with missing values in the **original** AOD tensor were first  
 446 filled with the spatial average of valid AOD data in each individual image. Then, the AOD tensor was  
 447 decomposed along each of three dimensions, while the dominant features in each dimension  
 448 determined by the corresponding rank values were applied to reconstruct the data tensor. By gradually  
 449 increasing the rank values and iteratively updating the initial filled values, the tensor can be  
 450 reconstructed to better delineate AOD distribution over space after several iterations.

451 To confirm the convergence, a small portion of observational AOD values were randomly held out  
 452 in advance, and the reconstructed values over these grids in each iteration were compared with these  
 453 hold-out data till the difference between them lower than 0.01 (a threshold to determine convergence,

Deleted: selected

Deleted:

Deleted: D

Deleted: were

Deleted: placed on

Deleted: greatly

Deleted: facilitated

Deleted: the

Deleted: of missing AOD information

Deleted: given the presence of prior knowledge

Deleted: More importantly, it significantly reduced the time required for convergence during the gap filling process.

Deleted: i.e.,

Field Code Changed

Deleted: , and Chen et al. (2014)

Deleted: .

469 a.k.a,  $\varepsilon_1$  in Table 2). Meanwhile, to make the computational burden manageable, the study region  
 470 (China in this study) was divided into 40 subregions (refer to Figure S2 for the spatial distribution of  
 471 these subregions), and the tensor completion was then performed over each individual region. Finally,  
 472 the reconstructed imagery were mosaiced to attain a national gap-free AOD map on each specific  
 473 date. During this step, a **smooth filter** was applied to solve the boundary effect when mosaicking two  
 474 adjacent maps. **Specifically**, data value on each overlapped grid at the boundary (50 km on the edge of  
 475 subregion) **was averaged via an inverse distance (the distance to the edge) weighting scheme**. In the  
 476 end, the mosaic AOD<sub>Terra</sub> image was retained as the final gap-free AOD product.

477 **Table 2.** The proposed tensor completion algorithm for AOD distribution reconstruction in AOD<sub>Terra</sub>.

**Input:** tensor  $\mathbf{A} \in \mathbf{R}^{N_1 \times N_2 \times N_3}$  with  $\Omega = \{(i, j, k): A_{ijk} \text{ is observed}\}$ , threshold  $T_1, T_2$   
**Output:** reconstructed entries  $\mathbf{A}' = \mathbf{A}^*(:, :, k^t) \in \mathbf{R}^{N_1 \times N_2}$

- 1: Initialize  $A_{ijk}^* = \begin{cases} A_{ijk} & (i, j, k) \in \Omega \\ \sum_i \sum_j A_{ijk} & (i, j, k) \notin \Omega \end{cases}$
- 2: **for**  $n_3 = N_3$  to 1 **do**
- 3:    $n_1 = n_2 = 0$
- 4:   **while**  $\varepsilon_1 > T_1$  **do**
- 5:      $n_1 = n_1 + 1, n_2 = n_2 + 1$
- 6:     Tucker Decomposition of  $\mathbf{A}^*$  with rank =  $\{n_1, n_2, n_3\}$ :  
        $\mathbf{A}^* = \mathbf{S} \times_1 \mathbf{U}^{(n_1)} \times_2 \mathbf{U}^{(n_2)} \times_3 \mathbf{U}^{(n_3)}$
- 7:      $\varepsilon_1 = \arg \min_{\Omega} \frac{1}{2} \|\mathbf{A} - \mathbf{A}^*\|^2$
- 8:      $\mathbf{A}_{\Omega}^* = \mathbf{A}_{\Omega}$
- 9:   **end while**
- 10:   **if**  $\arg \min_{\Omega} \frac{1}{2} \|\mathbf{A} - \mathbf{A}^*\|^2 < T_2$  **then**
- 11:     **break**;
- 12:   **end if**
- 13: **end for**

Deleted: weighted average

Deleted: method

Deleted: , i.e.,

Deleted: averaging the

Deleted: as

Deleted: weighted by

Deleted: the distance to the edge

### 478 3.3 PM<sub>x</sub> concentration estimation

479 In this study, the **widely used** random forest method was applied to establish regression models  
 480 for PM<sub>2.5</sub> and PM<sub>10</sub> concentration **estimation**. Ground measured PM<sub>2.5</sub> (or PM<sub>10</sub>) concentration data  
 481 were used as the learning target while **gap filled** AOD, aerosol components (AER<sub>comp</sub>), meteorological  
 482 factors (MET), digital elevation model (DEM), NDVI, land cover information (LC), and population  
 483 were used as regressors. The **random forest regression** model can be generally formulated as:

$$484 \quad \text{PM}_x = \text{RF}(AOD, AER_{comp}, MET, DEM, NDVI, POP, LC, month) \quad (1)$$

Formatted: Subscript

Deleted: mapping

Deleted: prediction

Deleted: f

495 where *month* is a categorical variable that was used to account for monthly varying relationships  
496 between AOD and  $PM_x$ . For validation,  $PM_{2.5}$  and  $PM_{10}$  measurements from 10% of monitoring sites  
497 were randomly held out to evaluate the predictive performance of each regression model. During the  
498 training process, 500 regression trees were used in each RF model, and each tree was grown on a  
499 bootstrap sample. The learning data set was randomly divided into two parts during the training process,  
500 with 80% used as the training set while the rest 20% for testing. In order to guarantee a larger value of  
501  $PM_{10}$  than  $PM_{2.5}$ ,  $PM_{2.5}$  estimations from Eq. (1) were used as one predictor in addition to factors used  
502 to predict  $PM_{2.5}$  when estimating  $PM_{10}$  concentration. Such a model can also significantly improve the  
503 prediction accuracy of  $PM_{10}$  given the prior  $PM_{2.5}$  information.

### 504 3.4 Point-surface data fusion

505 Ground measured  $PM_{2.5}$  and  $PM_{10}$  concentration data were further fused with their gridded  
506 estimations to enhance the data accuracy of  $PM_x$  data after 2014. Here, the well-known optimal  
507 interpolation (OI) method was applied to perform point-surface fusions between two different types  
508 datasets. Please refer to Bai et al. (2022) and Li et al. (2022) for a more detailed description of the OI  
509 method used to fuse  $PM_x$  concentration data. In this study, a modified scheme was developed to select  
510 neighboring observations. To avoid an isotropic interpolation effect, here we only used 30 ground  
511 observations with land cover, terrain and atmospheric conditions similar to those at the analyzed grid  
512 cell to estimate the innovation that should be assigned to the background value at the given grid. In  
513 other words, a similarity measure was first estimated between the analyzed grid cell and neighboring  
514 sites in terms of land cover, DEM, and atmospheric conditions. The 30 observations with similar  
515 background fields were then used in the OI procedure to correct possible bias in gridded  $PM_x$   
516 estimations. Such a treatment can help exclude those observations with different ambient background,  
517 e.g., one site not far from the given grid but separated by a high mountain, thereby avoiding the possible  
518 propagation of antiphase corrections to data over adjacent grids.

## 519 4 Results and discussion

### 520 4.1 Data accuracy of gap-free AOD in LGHAP

521 Table 3 summarizes the data accuracy of gap-free AOD dataset generated in this study. For  
522 comparison, the data accuracy of each original AOD dataset was also assessed. Since *in situ* AOD

Formatted: Font: Italic

Deleted: data

Formatted: Subscript

Deleted: cross

Deleted: validate

Deleted: PM

Formatted: Subscript

Deleted: 2021

Deleted: 2021

Formatted: Subscript

Formatted: Subscript

529 measurements were not used as data input when reconstructing missing AOD information, thereby the  
 530 gap-free AOD can be directly compared with *in situ* AOD measurements from AERONET. As  
 531 indicated, all these AOD datasets are in a good agreement with *in situ* AOD measurements. Generally,  
 532 AODs from MODIS onboard Terra and Aqua have an almost identical data accuracy, which is also  
 533 among the highest when comparing with other datasets (R=0.95 and RMSE=0.14). AODs from  
 534 AATSR show a comparable accuracy with that of MODIS, but with a relatively low correlation with  
 535 ground-based AOD measurements. AODs from MISR, POLDER and VIIRS exhibit a similar bias  
 536 level, with R varying from 0.80 to 0.92 and RMSE ranging from 0.22 to 0.29. In contrast, AOD<sub>M2</sub> data  
 537 have the poorest accuracy among these eight gridded AOD datasets (R=0.77 and RMSE=0.36), even  
 538 though AOD data from AERONET and satellite observations like MODIS had been already  
 539 assimilated. This indicates the presence of large biases in AOD<sub>M2</sub> and thus these AOD<sub>M2</sub> data cannot  
 540 be used solely to delineate AOD distributions over space.

Deleted: similar

541 **Table 3.** Data accuracy of original and gap-free AOD datasets used and/or generated in this study. The  
 542 expected error (EE) was defined as  $\pm 0.05 + 0.15 \times \text{AOD}_{\text{site}}$ .

Dataset	N	R	RMSE	MAE	Below EE (%)	Within EE (%)	Above EE (%)
Terra/MODIS	6731	0.95	0.13	0.07	8.94	78.73	12.33
Aqua/MODIS	6079	0.95	0.14	0.08	8.24	79.45	12.30
Terra/MISR	638	0.90	0.29	0.13	21.63	73.51	4.86
NPP/VIIRS	3839	0.80	0.22	0.16	7.03	44.93	48.03
Envisat/AATSR	434	0.92	0.11	0.07	17.74	73.96	8.29
PARASOL/POLDER	1733	0.92	0.24	0.17	5.14	40.22	54.65
MERRA-2	22067	0.77	0.36	0.20	32.97	51.76	15.27
LGHAP	24861	0.91	0.21	0.13	12.27	59.00	28.73

543  
 544 Compared to the first seven gridded AOD datasets, the LGHAP AOD dataset has an accuracy  
 545 slightly worse than the original MODIS AOD product but comparable to AODs from MISR, POLDER  
 546 and MERRA-2, with R of 0.91 and RMSE equaling to 0.21 compared to ground-based AOD  
 547 observations. Nevertheless, the gap-filled AOD appeared to overestimate ground-based AOD  
 548 observations, and this could be due to the involvement of AODs from VIIRS and POLDER as these  
 549 two products significantly overestimated ground AOD observations, which can be indicated by the

Deleted: Compared with AODs from MODIS

Deleted: caused by

Deleted: of

Deleted: significantly overestimated

555 proportion of data pairs above the expected error (EE). On the other hand, significant underestimations  
556 in AOD<sub>M2</sub> were not introduced to the LGHAP AOD as the former had a below EE ratio of 32.97%  
557 which was only 12.27% in the latter. These results indicate that the LGHAP AOD data are more likely  
558 to resemble AOD distributions revealed by satellite observations rather than AOD<sub>M2</sub>, endorsing the  
559 advantages of involving multisensory satellite AOD observations to support missing AOD  
560 reconstruction. Figure 2 further compares the data accuracy of original AOD<sub>Terra</sub> and the reconstructed  
561 data over different regions of China. It is indicative that the purely reconstructed data have an accuracy  
562 (R=0.88 and RMSE=0.26) lower than the original AOD<sub>Terra</sub> (R=0.95 and RMSE=0.13) across China,  
563 especially in South China where the reconstructed data were significantly underestimated the ground-  
564 based AOD observations. Possible reasons for this effect could be attributed to extensive data gaps in  
565 satellite AOD retrievals due to frequent and extensive cloud covers over there (refer to Figure S3 for  
566 the distribution of mean data integrity of AOD<sub>Terra</sub> during 2000–2020), and the scarce AOD  
567 observations significantly limit the learning capacity in space and temporal domain during the tensor  
568 completion process. In other words, limited observations in satellite imageries greatly reduced the  
569 learning performance from the sparse tensor. Even though, the purely reconstructed data exhibit a bias  
570 level comparable to AOD retrievals from several satellite instruments, e.g., MISR, VIIRS, and  
571 POLDER. This demonstrates the good performance of the proposed tensor completion method in  
572 reconstructing missing AOD information. By combining the reconstructed data with original AOD<sub>Terra</sub>,  
573 we obtained a 21-year-long gap free high-resolution (daily/1-km) AOD dataset with satisfying  
574 accuracy (R=0.91 and RMSE=0.21).

Deleted: such

Deleted: gap-free

Deleted: justifying

Deleted: ple

Deleted: help reconstruct

Deleted: values

Deleted: for the relatively poor accuracy of AOD reconstructions in this region

Deleted: over there

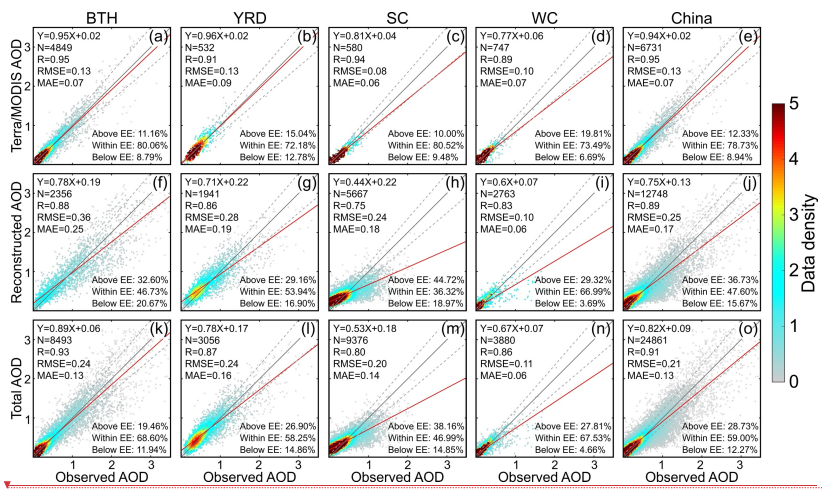
Deleted: which

Deleted: observations from Terra/MODIS

Deleted: spatially complete

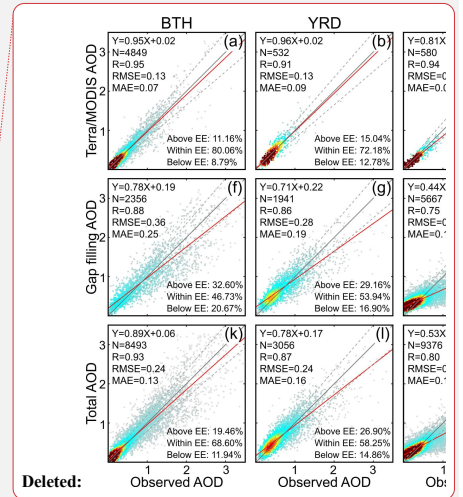
Deleted: product





**Figure 2.** Scatter plots between ground observed and satellite-based AOD data in different regions of China. (a–e) original Terra/MODIS AOD, (f–j) reconstructed AOD, and (k–o) combined AOD between original and reconstructed data. BTH, YRD, SC, and WC refers to regions of Beijing-Tianjin-Hebei, Yangtze River Delta, South China, and West China, respectively.

In Figure 3 we presented a comparison of AOD time series between the LGHAP dataset and ground observations at three AERONET sites under different air pollution levels. As shown, the AOD time series from LGHAP are temporally continuous whereas data gaps are common in AERONET observations. Generally, AODs from LGHAP are well reconstructed with respect to the temporal variations of aerosol loading at these three sites, with R ranging from 0.77 to 0.90 and RMSE varying between 0.11 and 0.21. For illustration, Figure 4 compares the spatial distribution of original and gap filled AOD on four days with different AOD<sub>Terra</sub> coverage over space. As shown, the missing AOD values were well reconstructed after gap filling, resembling a smooth and reasonable AOD distribution over space, even over regions with very limited prior AOD observations from Terra/MODIS (e.g., Figure 4d). As indicated in Figures 4a and 4c, the high AOD loading was also properly reconstructed even though no prior information was provided by AOD<sub>Terra</sub>. Since AERONET AOD observations were not used as a data input when generating the LGHAP AOD dataset, these independent validation results clearly demonstrated the high accuracy of the LGHAP AOD product as well as a good performance of the proposed AOD gap filling approach.



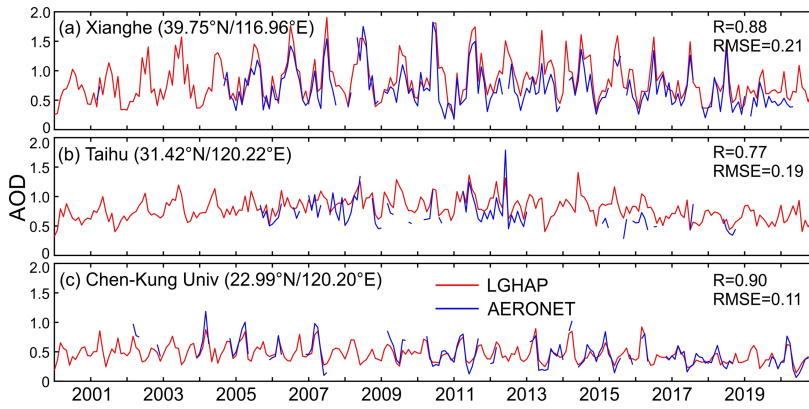
Deleted: only  
 Deleted: both  
 Deleted: combined

Deleted: reconstructed

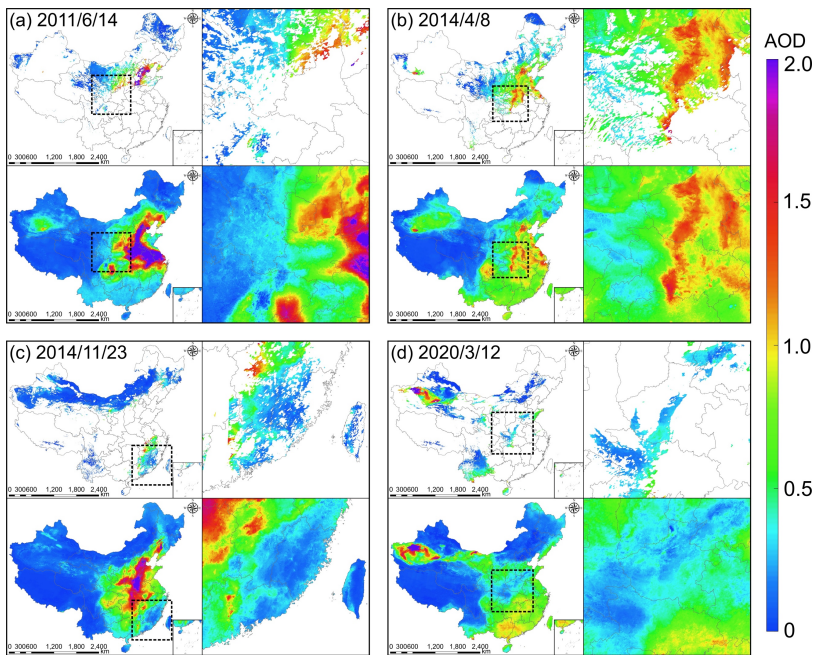
Deleted: ground-based

Deleted: full-coverage (gap free)

Deleted: mapping



615  
616 **Figure 3.** Comparison of monthly AOD time series from LGHAP and AERONET at three different  
617 stations in China. Latitude and longitude information of each site was given in brackets.



618  
619 **Figure 4.** Spatial patterns of the reconstructed AOD under different baseline AOD coverage ratios. In  
620 each sub-diagram, the upper panel presents the original AOD distribution from Terra/MODIS while

621 the gap-filled imagery is shown below. The zoom-in views of the outlined regions are shown in the  
622 right part.

Deleted: ¶

623 Since the final gap-free AOD product was generated mainly by integrating a set of data tensor  
624 of gridded AOD with AOD estimations *from in situ air quality measurements*, the relative contribution  
625 of each product to the final gap-free dataset is worth being investigated. In this study, a data coverage  
626 ratio weighted nonlinear correlation coefficient was proposed to examine the relative contribution of  
627 each gridded product to the LGHAP AOD dataset. The nonlinear correlation coefficient was used to  
628 assess the mutual information between two variables (Sun et al., 2021; Wang et al., 2005), while the  
629 data coverage ratio was multiplied to indicate the overall contribution of one product to the final fused  
630 dataset (refer to Text S2 for the definition of this indicator). As shown in Figure 5, the relative  
631 contribution of each gridded product varied with time and the input data sources. In the early two years  
632 (2000–2001), the AOD distribution in gap-free imageries was determined largely by AOD<sub>Terra</sub> (81%),  
633 whereas this ratio decreased to about 30% when many other products were involved, especially AOD  
634 from Aqua and PARASOL. With the advent of VIIRS and the loss of PARASOL after 2012, the  
635 relative contribution changed drastically as AOD from MODIS and VIIRS played the dominant roles  
636 in reconstructing AOD distribution. Note the relative contribution of AOD<sub>M2</sub> remained lower than 10%,  
637 indicative of the greater importance of satellite observations in generating the LGHAP AOD product.

Deleted: *in situ*

Formatted: Font: Italic

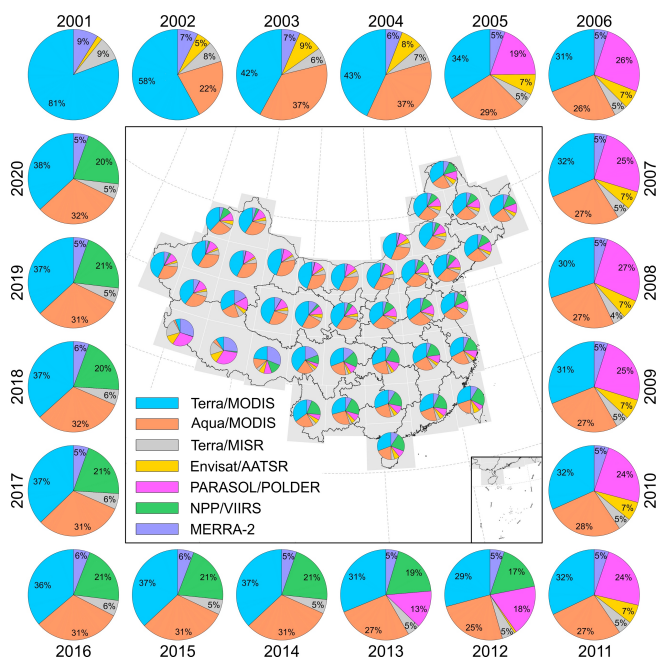
Deleted:

638 With respect to the temporally averaged contribution in each subregion, it shows that the  
639 relative contribution of each product also varied significantly across regions. Generally, AOD from  
640 MODIS aboard Terra and Aqua played the most important role (>60%) in generating the LGHAP  
641 AOD product, except over the southwest part of the country (Tibet plateau) where AOD<sub>M2</sub> contributed  
642 most. This is *largely associated with the fact that* data gaps are abnormally high in satellite observations  
643 over this region because of the vast and long-lasting snow cover (refer to Figure S3 for the data  
644 integrity distribution). Consequently, AOD<sub>M2</sub> would play an important role in reconstructing AOD  
645 distribution over such regions. *Note that the relative contribution of AOD estimations from in situ air*  
646 *quality measurements were not accounted for in the current analysis because of incomparable spatial*  
647 *coverage of in situ data contrast to gridded AOD products, and this does not imply the contribution of*  
648 *in situ AOD estimations being negligible.* Overall, the results shown here clearly highlight the success  
649 of big data analytics in generating the LGHAP AOD dataset via integrative efforts from diversified  
650 data sources.

Deleted: reasonable since

Formatted: Font: Italic

Formatted: Font: Italic



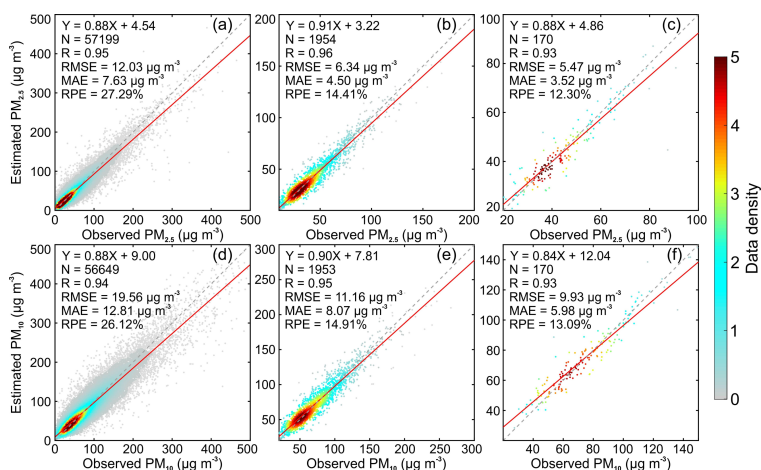
655  
 656 **Figure 5.** Spatiotemporal variations of the relative contribution of each gridded AOD product to the  
 657 generation of LGHAP AOD dataset. The relative contribution was estimated as the data coverage ratio  
 658 weighted nonlinear correlation coefficient (please refer to Text S2 in the supplementary information  
 659 for the arithmetic theory to calculate this measure). The annual mean shown outside is the national  
 660 averaged contribution in each individual year while the regional mean shown on the map was averaged  
 661 over the past 21-year in each subregion.

662 **4.2 Data accuracy of PM<sub>2.5</sub> and PM<sub>10</sub> estimations**

663 By taking advantage of the gap-filled AOD, daily 1-km resolution PM<sub>2.5</sub> and PM<sub>10</sub> concentration  
 664 data in China were [then](#) estimated via an ensemble learning approach. Figure S4 shows the sample-  
 665 based cross validation accuracy of two prediction models. It shows that the original daily PM<sub>2.5</sub>  
 666 prediction model had a sample-based cross validation R<sup>2</sup> of 0.79 and RMSE of 20.04 μg m<sup>-3</sup>. This  
 667 accuracy is comparable to our previous study (Bai et al., 2019a), but slightly worse than those reported  
 668 in some recent studies (Table 4). In contrast, PM<sub>10</sub> had a much higher prediction accuracy, with R<sup>2</sup> of  
 669 0.90 and RMSE of 21.06 μg m<sup>-3</sup> for the daily product. This good performance should be attributed to

670 the involvement of  $PM_{2.5}$  estimations as a predictor in the  $PM_{10}$  prediction model. Figure 6 shows the  
 671 site-specific (held-out in advance) validation accuracy of daily, monthly, and annual mean  $PM_{2.5}$  and  
 672  $PM_{10}$  concentration in LGHAP. As shown, the site-specific validation results indicated that the final  
 673 full-coverage (gap free) daily  $PM_{2.5}$  and  $PM_{10}$  concentration data are in a good agreement with ground-  
 674 based measurements, with R of 0.95 and RMSE of  $12.03 \mu g m^{-3}$  for  $PM_{2.5}$  while R of 0.94 and RMSE  
 675 of  $19.56 \mu g m^{-3}$  for  $PM_{10}$ . Overall,  $PM_x$  data in LGHAP are not only spatially complete with a finer  
 676 resolution but have a comparable accuracy with previous studies.

Formatted: Subscript



677  
 678 **Figure 6.** Scatter plots between observed and estimated  $PM_{2.5}$  and  $PM_{10}$  concentration. (a–c)  
 679 respectively denotes daily, monthly, and annual mean  $PM_{2.5}$  validation results, while (d–f) are for  $PM_{10}$   
 680 concentration. The ground measurements were acquired from 30 independent air quality monitoring  
 681 sites that were randomly held-out before the model training.

682 **Table 4.** Comparison of the data quality of  $PM_{2.5}$  from LGHAP with other related studies.

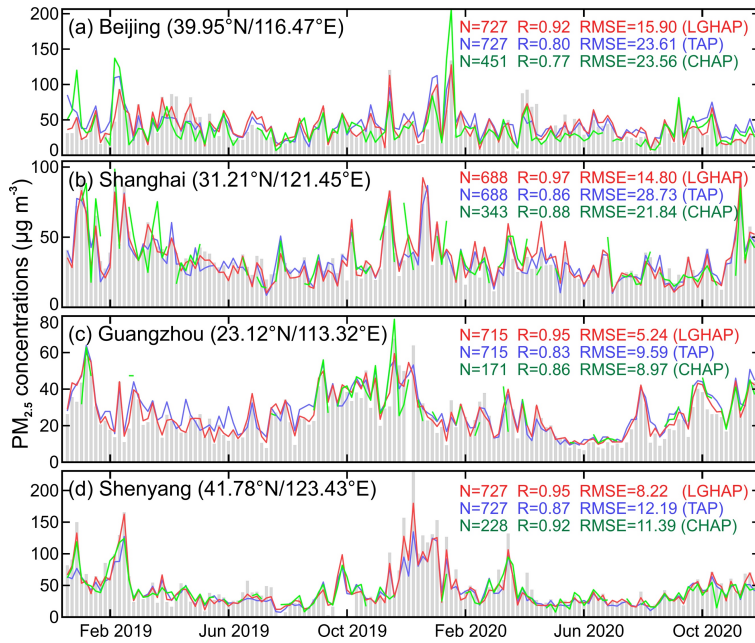
Source	Gap-free	Resolution	Time range	$R^2$	RMSE ( $\mu g m^{-3}$ )
Wei et al. (2021a)	No	1 km	2000~2018	0.86~0.90	10.09~18.39
Geng et al. (2021)	Yes	10 km	2000~2021	0.80~0.88	13.90~22.10
Xue et al. (2019)	Yes	10 km	2000~2016	0.61	27.80
Chen et al. (2018)	No	10 km	2005~2016	0.83	28.10
Lyu et al. (2019)	Yes	12 km	2014~2017	0.64	24.80
Ma et al. (2016)	No	10 km	2004~2013	0.79	27.42

Huang et al. (2021)	No	1 km	2013~2019	0.88	15.73
Xiao et al. (2018)	Yes	10 km	2013~2017	0.79	21.00
LGHAP PM <sub>2.5</sub>	Yes	1 km	2000~2020	0.90	12.03

683

684 Figure 7 presents a two-year-long comparison of PM<sub>2.5</sub> concentration time series from LGHAP  
685 and two other open access datasets with PM<sub>2.5</sub> measurements sampled at four United States Embassy  
686 in China. Since this ground-based dataset has been seldomly noticed and used, it can be applied as an  
687 independent dataset to fairly evaluate the accuracy of these three machine-learned PM<sub>2.5</sub> estimations.  
688 As shown, all these three datasets well reconstructed temporal variations of PM<sub>2.5</sub> from 2019 to 2020.  
689 Temporally, LGHAP and TAP are continuous while CHAP suffers from significant data gaps because  
690 no gap filling was applied when generating the dataset. Compared with the other two datasets, LGHAP  
691 PM<sub>2.5</sub> data had a better agreement with ground-based PM<sub>2.5</sub> measurements. This high accuracy could  
692 be partially due to the fusion of *in situ* PM<sub>2.5</sub> data measured at adjacent sites via the OI method. Figure  
693 S5 compares PM<sub>2.5</sub> time series from LGHAP with PM<sub>2.5</sub> measurements sampled at five United States  
694 Embassy in China. It is indicative that historical PM<sub>2.5</sub> variations over these five cities were well  
695 reconstructed in LGHAP, even over years before 2014 at which PM<sub>2.5</sub> measurements from state-  
696 control monitoring sites were not available. Note PM<sub>2.5</sub> estimations appeared to significantly  
697 underestimate PM<sub>2.5</sub> concentration sampled at the Embassy in Beijing before 2013. Considering the  
698 reconstructed AOD time series agreed well with AERONET AOD in Beijing (Figure 3a), and the  
699 model performed well in predicting historical PM<sub>2.5</sub> in Shanghai during the synchronous time period  
700 (Figure S5b), we are more willing to attribute this issue to significant PM<sub>2.5</sub> overestimations by the US  
701 Embassy during that period. Overall, these independent validation results collectively indicate a good  
702 accuracy of PM<sub>2.5</sub> in LGHAP dataset.

Deleted: method

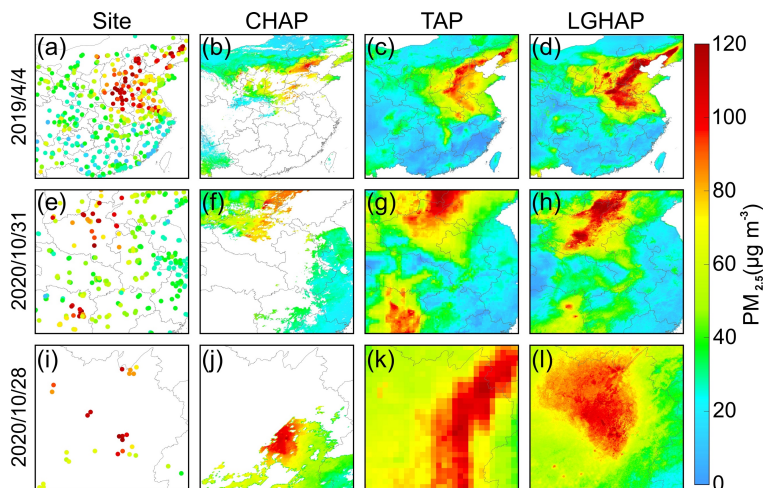


704 **Figure 7.** Comparison of PM<sub>2.5</sub> concentration time series between LGHAP (red line) and two open  
 705 datasets (blue: TAP, green: CHAP). Here, hourly PM<sub>2.5</sub> concentrations measured by four United States  
 706 Embassy in China from 2019 to 2020 (grey bar) were used as an independent PM<sub>2.5</sub> dataset to validate  
 707 these three daily products. CHAP and TAP are two open access datasets providing PM<sub>2.5</sub>  
 708 concentration that were created by Wei et al. (2021a) and Geng et al. (2021) respectively.  
 709

710  
 711 In Figure 8 we compared the spatial distribution of PM<sub>2.5</sub> that was reconstructed by different  
 712 datasets. Compared to LGHAP and TAP, PM<sub>2.5</sub> data from CHAP are not gap free since the spatial  
 713 coverage is determined by the AOD data coverage in the MAIAC product. Compared to TAP, LGHAP  
 714 PM<sub>2.5</sub> data have a finer resolution (1 km versus 10 km), enabling us to examine PM<sub>2.5</sub> variations in  
 715 space with more details. Overall, LGHAP has a better performance in reconstructing PM<sub>2.5</sub> spatial  
 716 distributions than the other two datasets. Reasons could be attributed to the following two aspects.  
 717 Firstly, *in situ* PM<sub>2.5</sub> measurements were fused with gridded PM<sub>2.5</sub> estimations using the OI method  
 718 when generating the final PM<sub>2.5</sub> product in LGHAP. This can help correct modeling biases in original  
 719 PM<sub>2.5</sub> estimations. Secondly, a set of satellite-based AOD retrievals were incorporated when

Formatted: Font: Italic

720 generating the full-coverage AOD product, which greatly helps reduce large biases in numerical AOD  
 721 simulations, yielding more accurate  $PM_{2.5}$  estimations in turn. This also highlights the great advantages  
 722 of using big data analytics methods to advance air pollution assessment.

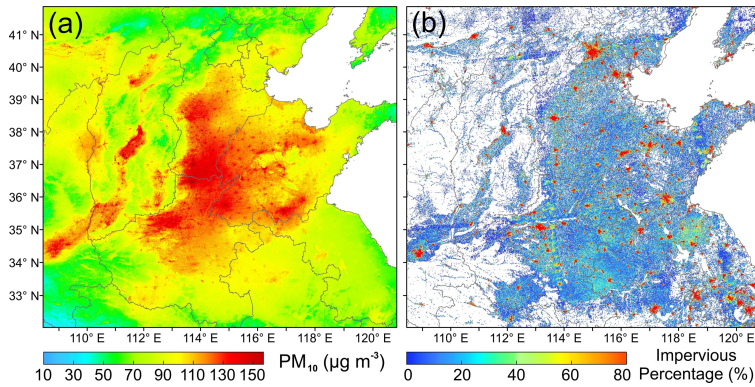


723  
 724 **Figure 8.** Comparison of  $PM_{2.5}$  distribution reconstructed by different  $PM_{2.5}$  concentration datasets.  
 725 From the left to right, it shows in situ  $PM_{2.5}$  concentration measurements, CHAP, TAP, and LGHAP,  
 726 respectively.

727  
 728 To illustrate the fine resolution of LGHAP dataset, we compared the annual mean  $PM_{10}$   
 729 concentration in 2019 with the proportion of impervious surface that was derived from 30-m resolution  
 730 land cover data in eastern China. As shown in Figure 9, the finer resolution of LGHAP dataset enables  
 731 us to easily recognize the “hot spot” regions with high  $PM_{10}$  loading. By referring to the impervious  
 732 surface distribution on the right, we found that these hot spots are mainly over cities and towns,  
 733 indicative of the presence of pollution island in urban regions. Owing to the involvement of such high-  
 734 resolution datasets, the spatial details of  $PM_{2.5}$  and  $PM_{10}$  can be then well recognized in LGHAP. The  
 735 finer spatial resolution advantage of the LGHAP dataset can be also demonstrated by comparisons of  
 736 spatial distribution of annual mean  $PM_{2.5}$  concentration that was revealed by four different datasets  
 737 shown in Figure S6.

Deleted: set





739  
740 **Figure 9.** Comparison of annual mean PM<sub>10</sub> concentration with the proportion of areas covered by  
741 impervious surface in eastern China.

742 **4.3 Long-term trends of haze pollution in China from 2000 to 2020**

743 The aerosol pollution trends in China can be better examined by taking advantage of LGHAP  
744 dataset given long temporal coverage, gap free and high-resolution superiorities. Severe haze  
745 pollutions such as PM<sub>2.5</sub> are oftentimes observed during winter half year (September–February). In  
746 this study, we first calculated mean PM<sub>2.5</sub> concentration in China during winter half year from 2000 to  
747 2020. As shown in Figure 10, severe haze pollution events were mainly observed in North China during  
748 the wintertime, especially over the adjacent region in Hebei-Shandong-Henan provinces. In addition,  
749 Sichuan basin and Fenwei plain also suffered from severe haze pollution. Temporally, severe haze  
750 pollution events occurred mainly from the late 2002 to early 2017, which were significantly reduced  
751 after 2017. Similar pattern can be also inferred from PM<sub>10</sub> concentration distributions shown in Figure  
752 S7.

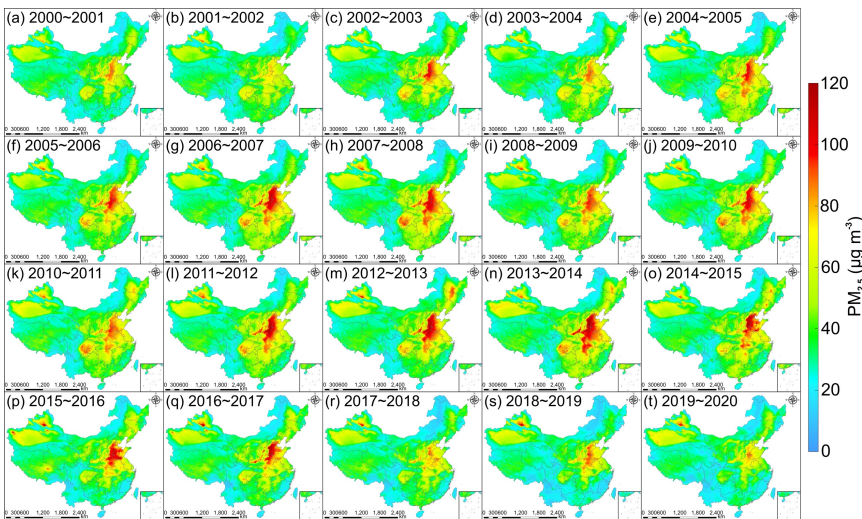
753 Figure 11 shows the temporal variations of the proportion of land areas covered by PM<sub>2.5</sub>  
754 concentration exceeding 35 µg m<sup>-3</sup> (the national ambient air quality standard for 24-hour PM<sub>2.5</sub>  
755 concentration given in GB 3095-2012). As shown in Figure 11a, severe PM<sub>2.5</sub> pollution occurred  
756 mainly during the wintertime in China, as more than one-third land areas (indicated by the blue lines)  
757 were exposed to unhealthy PM<sub>2.5</sub> pollutants. Meanwhile, an apparent inflection was observed in 2007,  
758 after which the number of episode days decreased drastically at more than one-third land area covered  
759 by PM<sub>2.5</sub> concentration exceeding 35 µg m<sup>-3</sup>. According to the proportion of land area covered with

- Deleted: the
- Deleted: wintertime
- Deleted: to
- Deleted:
- Deleted: wintertime
- Deleted: wintertime
- Deleted: during the wintertime

- Deleted: hazardous

768 annual mean  $PM_{2.5}$  concentration greater than  $35 \mu g m^{-3}$ , the variation of haze pollution in China can  
 769 be generally divided into three different periods during the past two-decades (Figure 11b). As indicated,  
 770 an increasing trend was observed from 2000 to 2007, during which land areas covered by  $PM_{2.5}$   
 771 concentration greater than  $35 \mu g m^{-3}$  had increased to near 40% at a pace of  $1.04\% a^{-1}$ . The second  
 772 period was from 2008 to 2013, during which the land area coverage ratio decreased at a rate of  $-0.21\%$   
 773  $a^{-1}$ . The third period started from 2014, after which the land area covered with  $PM_{2.5}$  concentration  
 774 more than  $35 \mu g m^{-3}$  had decreased drastically, at a **pace** of  $-2.23\% a^{-1}$ .

Deleted: rate



775  
 776 **Figure 10.** Spatial distribution of **mean**  $PM_{2.5}$  concentration from LGHAP during **winter half year**  
 777 **(September–February)** from 2000 to 2020 in China.

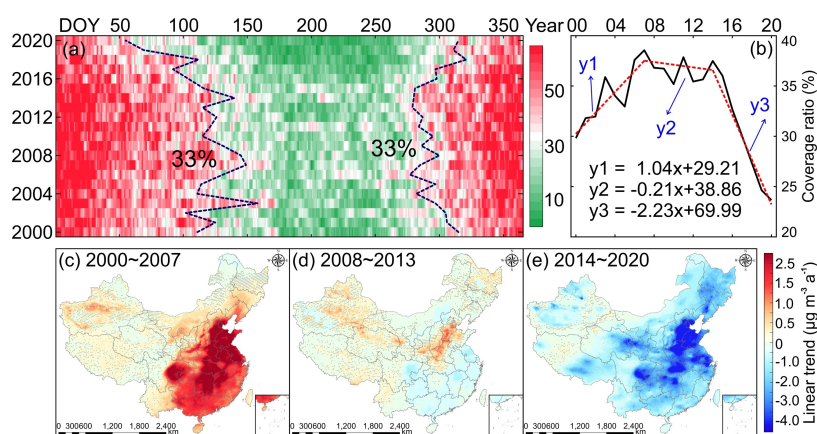
Deleted: wintertime (September to February) averaged

779 Figure 11c–e presents the linear trend of  $PM_{2.5}$  concentration during these three specific periods,  
 780 from which we observed that significant  $PM_{2.5}$  variations occurred mainly over eastern part of the  
 781 country where resides two-thirds of the population. A near ubiquitous  $PM_{2.5}$  increasing trend was  
 782 observed during 2000–2007, with significant increase ( $>1.0 \mu g m^{-3} a^{-1}$ ) mainly observed in eastern  
 783 China. During the second period,  $PM_{2.5}$  concentration over most regions shows a small decreasing  
 784 trend except in the Ji-Lu-Yu region where an increasing trend was still observed. Apparent decreasing  
 785 trend was observed over most parts of the country after 2014, indicative of significant reductions in  
 786  $PM_{2.5}$  loading across China. This trend distribution is in line with our previous **finding** that was derived

Deleted: s

Deleted: study

791 using the annual mean  $PM_{2.5}$  concentration dataset generated by the Dalhousie University (Bai et al.,  
 792 2019b). However, differences were still observed in terms of the regions where significant decreasing  
 793 trends were present. Most significant decreasing trends were mainly observed in Sichuan basin and  
 794 Pearl River Delta in the previous study. However, regions with drastic  $PM_{2.5}$  decrease were found  
 795 mainly in the North China where severe haze pollution events were oftentimes reported. Similar  
 796 variation patterns can be also inferred from  $PM_{10}$  (Figure S8) and AOD (Figure S9). Overall, the  
 797 LGHAP dataset provides us a gridded perspective to better examine long-term variations of haze  
 798 pollution in China during the past two decades.

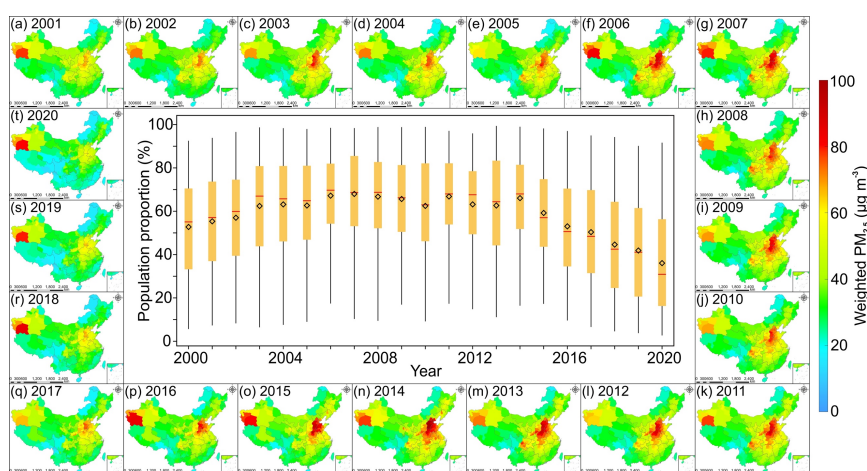


799  
 800 **Figure 11.** Temporal variations of the proportion of land areas covered with  $PM_{2.5}$  concentration  
 801 exceeding  $35 \mu\text{g m}^{-3}$  and  $PM_{2.5}$  trends during three different periods. (a) Temporal variations of the  
 802 land coverage ratio with daily  $PM_{2.5}$  concentration exceeding  $35 \mu\text{g m}^{-3}$  from 2000 to 2000. (b) same  
 803 as (a) but for annual mean  $PM_{2.5}$  concentration. (c–e)  $PM_{2.5}$  trends during periods of 2000–2007,  
 804 2008–2013, and 2014–2020. The dotted regions imply trend estimations are statistically insignificant at the  
 805 95% confidence interval.

806 **4.4 Population exposure to  $PM_{2.5}$  pollution in China**

807 By taking advantage of fine resolution LGHAP  $PM_{2.5}$  concentration and gridded population data,  
 808 population exposure to  $PM_{2.5}$  pollution across China over the past two decades were estimated. Figure  
 809 12 shows the spatial distribution of population weighted  $PM_{2.5}$  concentration and the proportion of  
 810 population exposed to  $PM_{2.5}$  concentration greater than  $35 \mu\text{g m}^{-3}$ . As shown, spatial distribution of

811 population weighted  $PM_{2.5}$  concentration resembles the spatial pattern of annual mean  $PM_{2.5}$   
812 concentration, with high values observed mainly in eastern and central China as well as northwest  
813 China. Nonetheless,  $PM_{2.5}$  sources in these two areas could be different. In northwest China, natural  
814 emissions could be the dominant source since very limited population resides there. In contrast, most  
815 population lives in eastern and central China with highly developed economy, and anthropogenic  
816 emissions thus might play more important roles in  $PM_{2.5}$  formation (Xin et al., 2015; Yang et al., 2011).  
817 In regard to the proportion of population exposed to the ambient with  $PM_{2.5}$  concentration greater than  
818  $35 \mu g m^{-3}$ , we observed that the annual mean population ratio exposure to unhealthy  $PM_{2.5}$  increased  
819 gradually from 50.60% in 2000 to 65.72% in 2007. During 2007–2014, the ratio varied with small  
820 changes (<5%), whereas a drastic decline was observed after 2014, with the annual mean proportion  
821 of population exposed to unhealthy  $PM_{2.5}$  was reduced from 63.81% in 2014 to 34.03% in 2020, even  
822 though the total population was increased from 1.37 billion to 1.41 billion during the synchronous  
823 period. Nonetheless, more than one-third population was still exposed to unhealthy  $PM_{2.5}$ , highlighting  
824 the requirement of further emission reduction actions to manage haze pollutions in China.



825  
826 **Figure 12.** Spatial distribution of population weighted  $PM_{2.5}$  concentration and the proportion of  
827 population exposed to  $PM_{2.5}$  concentration greater than  $35 \mu g m^{-3}$ . Annual and daily LGHAP  $PM_{2.5}$   
828 concentration data were used for the calculation of weighted  $PM_{2.5}$  and the proportion of population  
829 exposure, respectively. The diamond and red line indicate the annual mean and median population  
830 proportion, respectively.

831 **5 Data availability**

832 The LGHAP dataset, consisting of gap free AOD, PM<sub>2.5</sub>, and PM<sub>10</sub> concentration with daily 1-  
833 km resolution from 2000 to 2020, are all publicly accessible. ~~All data were~~ provided in the NetCDF  
834 format and data in each individual year were archived in a zip file. For AOD, the dataset has a disk  
835 storage size of near 27 GB in total, which ~~is available~~ at <https://doi.org/10.5281/zenodo.5652257> (Bai et  
836 al., 2021a). PM<sub>2.5</sub> (38 GB) and PM<sub>10</sub> (48 GB) concentration data can be acquired from  
837 <https://doi.org/10.5281/zenodo.5652265> (Bai et al., 2021b) and <https://doi.org/10.5281/zenodo.5652263> (Bai  
838 et al., 2021c), respectively. Additionally, monthly and annual mean datasets were also provided, which  
839 ~~is publicly available at~~ <https://doi.org/10.5281/zenodo.5655797> (Bai et al., 2021d) and  
840 <https://doi.org/10.5281/zenodo.5655807> (Bai et al., 2021e), respectively. ~~In addition to these datasets,~~  
841 Python, Matlab, R, and IDL codes that can be used to read and visualize these data were provided as  
842 well.

843 **6 Conclusion**

844 In this study, a big data analytics method was developed for generating a LGHAP dataset to  
845 advance research in earth system science and environment management. With integrative efforts of  
846 fusing AOD features extracted from a set of AOD data tensors and knowledge transfer in statistical  
847 data mining from diverse air quality indicators, a LGHAP aerosol dataset providing 21-year-long  
848 (2000–2020) gap-free AOD, PM<sub>2.5</sub>, and PM<sub>10</sub> concentration data with daily 1-km resolution in China,  
849 was generated. Gap-filled AOD imageries were firstly generated by reconstructing AOD distribution  
850 in AOD<sub>Terra</sub> via *synergistically* fusing AOD features recognized from diversified satellites and  
851 numerical models as well as *in situ* data through tensor completion. Compared to ground-based AOD  
852 measurements, the gap-filled AOD data exhibit a satisfying prediction accuracy and good performance  
853 in delineating AOD variations over space and time. To our knowledge, this is the first thrust of  
854 generating long-term high-resolution AOD dataset with gap free nature in China.

855 PM<sub>2.5</sub> and PM<sub>10</sub> concentration data were then estimated using an ensemble learning approach by  
856 taking advantage of the generated gap-free AOD imageries. Ground validation results also indicate  
857 good accuracies of these two gridded products, showing a comparable bias level with many previous  
858 studies. Compared with other open access daily PM<sub>2.5</sub> concentration datasets, the LGHAP PM<sub>2.5</sub>  
859 dataset performs well due to the vantage of having gap free and fine resolution products. With this gap

Deleted: The daily map

Deleted: was

Deleted: can be found

Deleted: can be acquired from

Moved (insertion) [1]

Deleted: LGHAP

Deleted: Global scale dataset is on the track and will be released to the public soon.

Formatted: Font: Italic

867 free and high-resolution dataset, the long-term variation trend of haze pollution in China over the past  
868 two decades was examined, and apparent inflections were observed in 2007 and 2014, at which PM<sub>2.5</sub>  
869 concentration was found to turn from an increasing path to decreasing in 2007 with a more drastic  
870 decline observed starting from 2014. Moreover, the LGHAP dataset provides us a gridded perspective  
871 to assess two-decade long population exposure to PM<sub>2.5</sub> pollution in China. In spite of a drastic decline  
872 in population exposure, there are still more than one-third population exposed to unhealthy PM<sub>2.5</sub>  
873 pollutants, highlighting the requirement of long-lasting actions to continue PM<sub>2.5</sub> related emission  
874 reduction.

875 Overall, these three gridded LGHAP aerosol products provide a long-term perspective on aerosol  
876 changes over different regions of China, and users are encouraged to use the LGHAP dataset to assess  
877 aerosol impacts on public health, air quality, climate, and ecosystem. The dataset has been publicly  
878 released online and is freely accessible via the links provided [in Section 5](#). Global scale dataset is on  
879 the track and will be released to the public soon.

Deleted: above

Moved up [1]: In addition to the LGHAP dataset, Python, Matlab, R, and IDL codes that can be used to read and visualize these data were provided as well. Global scale dataset is on the track and will be released to the public soon.

#### 880 **Author contributions**

881 The study was completed with cooperation between all authors. KB, KL, JG, ZL and N.B.C conceived  
882 of the idea behind generating the LGHAP dataset. KL, KB, and ZT developed the method and KB  
883 wrote the paper. KL, KB, K.T.L, and MM conducted the data analyses. JG and ZL provided  
884 atmospheric visibility and in situ AOD data, respectively. All authors discussed the results and  
885 proofread the paper.

#### 886 **Competing interests**

887 The authors declare that they have no conflict of interest.

#### 888 **Acknowledgments**

889 The authors are grateful to [the editor and two anonymous referees for their constructive comments and](#)  
890 [suggestions in improving this manuscript. Also, the authors would like to thank](#) all organizations and  
891 groups for providing essential datasets that were used in this study. The MAIAC AOD was acquired  
892 from <https://lpdaac.usgs.gov/products/mcd19a2v006/>. The MISR AOD was acquired from  
893 <https://asdc.larc.nasa.gov/project/MISR>. The VIIRS AOD was acquired from

899 <https://earthdata.nasa.gov/earth-observation-data/near-real-time/download-nrt-data/viirs-nrt>. The  
900 AATSR AOD was acquired from <https://climate.esa.int/en/projects/aerosol/data/>. The POLDER AOD  
901 was acquired from <https://www.grasp-open.com/products/polder-data-release/>. The aerosol  
902 diagnostics including AOD and aerosol components from MERRA-2 were acquired from  
903 [https://disc.gsfc.nasa.gov/datasets/M2T1NXAER\\_5.12.4/summary?keywords=MERRA2](https://disc.gsfc.nasa.gov/datasets/M2T1NXAER_5.12.4/summary?keywords=MERRA2). AOD from  
904 AERONET was acquired from [https://aeronet.gsfc.nasa.gov/new\\_web/aerosols.html](https://aeronet.gsfc.nasa.gov/new_web/aerosols.html). Meteorological  
905 factors were retrieved from the latest ERA-5 reanalysis and can be reached at  
906 <https://cds.climate.copernicus.eu/cdsapp#!/dataset/reanalysis-era5-single-levels?tab=overview>.  
907 Atmospheric visibility data were acquired from the national meteorological information center at  
908 <http://data.cma.cn/en>. Ground-based air pollutants concentration was acquired from  
909 <https://air.cnemc.cn:18007/>. Gridded Population data were acquired from <https://www.worldpop.org/>  
910 while DEM was acquired from <https://www.resdc.cn/>. Monthly NDVI data were acquired from  
911 <https://lpdaac.usgs.gov/products/mod13a3v061/>. Land cover data were acquired from  
912 <http://www.globallandcover.com/defaults.html?src=/Scripts/map/defaults/browse.html&head=brows>  
913 [e&type=data](http://www.globallandcover.com/defaults.html?src=/Scripts/map/defaults/browse.html&head=browse&type=data) and <https://zenodo.org/record/4417810#.YSxD844zYuW>.

#### 914 **Financial support**

915 This study was supported by the National Natural Science Foundation of China (grants 42171309 and  
916 41701413), and the Shanghai Committee of Science and Technology (grant 20ZR1415900).

Deleted: ¶

918 **References**

- 919 Bai, K., Chang, N.-B. and Chen, C.-F.: Spectral Information Adaptation and Synthesis Scheme  
920 for Merging Cross-Mission Ocean Color Reflectance Observations From MODIS and VIIRS, IEEE  
921 Trans. Geosci. Remote Sens., 54(1), 311–329, doi:10.1109/TGRS.2015.2456906, 2016.
- 922 Bai, K., Li, K., Chang, N.-B. and Gao, W.: Advancing the prediction accuracy of satellite-based  
923 PM<sub>2.5</sub> concentration mapping: A perspective of data mining through in situ PM<sub>2.5</sub> measurements,  
924 Environ. Pollut., 254, 113047, doi:10.1016/j.envpol.2019.113047, 2019a.
- 925 Bai, K., Ma, M., Chang, N.-B. and Gao, W.: Spatiotemporal trend analysis for fine particulate  
926 matter concentrations in China using high-resolution satellite-derived and ground-measured PM<sub>2.5</sub>  
927 data, J. Environ. Manage., 233, 530–542, doi:10.1016/j.jenvman.2018.12.071, 2019b.
- 928 Bai, K., Li, K., Wu, C., Chang, N.-B. and Guo, J.: A homogenized daily in situ PM<sub>2.5</sub>  
929 concentration dataset from the national air quality monitoring network in China, Earth Syst. Sci. Data,  
930 12(4), 3067–3080, doi:10.5194/essd-12-3067-2020, 2020a.
- 931 Bai, K., Li, K., Guo, J., Yang, Y. and Chang, N.-B.: Filling the gaps of in situ hourly PM<sub>2.5</sub>  
932 concentration data with the aid of empirical orthogonal function analysis constrained by diurnal cycles,  
933 Atmos. Meas. Tech., 13(3), 1213–1226, doi:10.5194/amt-13-1213-2020, 2020b.
- 934 Bai, K., Li, K., Guo, J. and Chang, N.-B.: Multiscale and multisource data fusion for full-coverage  
935 PM<sub>2.5</sub> concentration mapping: Can spatial pattern recognition come with modeling accuracy?, ISPRS  
936 J. Photogramm. Remote Sens., 184, 31–44, doi: 10.1016/j.isprsjprs.2021.12.002, 2022.
- 937 Bai, K., Li, K., Tan, Z., Han, D., and Guo, J.: Daily 1-km gap-free AOD grids in China, v1 (2000–  
938 2020) [data set], <https://doi.org/10.5281/zenodo.5652257>, 2021a.
- 939 Bai, K., Li, K., Tan, Z., Han, D., and Guo, J.: Daily 1-km gap-free PM<sub>2.5</sub> grids in China, v1 (2000–  
940 2020) [data set], <https://doi.org/10.5281/zenodo.5652265>, 2021b.
- 941 Bai, K., Li, K., Tan, Z., Han, D., and Guo, J.: Daily 1-km gap-free PM<sub>10</sub> grids in China, v1 (2000–  
942 2020) [data set], <https://doi.org/10.5281/zenodo.5652263>, 2021c.
- 943 Bai, K., Li, K., Tan, Z., Han, D., and Guo, J.: Monthly averaged 1-km gap-free AOD, PM<sub>2.5</sub> and  
944 PM<sub>10</sub> grids in China, v1 (2000–2020) [data set], <https://doi.org/10.5281/zenodo.5655797>, 2021d.
- 945 Bai, K., Li, K., Tan, Z., Han, D., and Guo, J.: Annual mean 1-km gap-free AOD, PM<sub>2.5</sub> and PM<sub>10</sub>  
946 grids in China, v1 (2000–2020) [data set], <https://doi.org/10.5281/zenodo.5655807>, 2021e.
- 947 Beckers, J. M. and Rixen, M.: EOF Calculations and Data Filling from Incomplete Oceanographic

Deleted: 2021. in revision



949 Datasets, *J. Atmos. Ocean. Technol.*, 20(12), 1839–1856, doi:10.1175/1520-  
950 0426(2003)020<1839:ECADFF>2.0.CO;2, 2003.

951 Bi, J., Belle, J. H., Wang, Y., Lyapustin, A. I., Wildani, A. and Liu, Y.: Impacts of snow and cloud  
952 covers on satellite-derived PM<sub>2.5</sub> levels, *Remote Sens. Environ.*, 221(October), 665–674,  
953 doi:10.1016/j.rse.2018.12.002, 2018.

954 Chang, N.-B., Bai, K. and Chen, C.-F.: Smart Information Reconstruction via Time-Space-  
955 Spectrum Continuum for Cloud Removal in Satellite Images, *IEEE J. Sel. Top. Appl. Earth Obs.*  
956 *Remote Sens.*, 8(5), 1898–1912, doi:10.1109/JSTARS.2015.2400636, 2015.

957 Che, H., Yang, L., Liu, C., Xia, X., Wang, Y., Wang, H., Wang, H., Lu, X. and Zhang, X.: Long-  
958 term validation of MODIS C6 and C6.1 Dark Target aerosol products over China using CARSNET  
959 and AERONET, *Chemosphere*, 236, 124268, doi:10.1016/j.chemosphere.2019.06.238, 2019.

960 Chen, G., Li, S., Knibbs, L. D., Hamm, N. A. S., Cao, W., Li, T., Guo, J., Ren, H., Abramson, M.  
961 J. and Guo, Y.: A machine learning method to estimate PM<sub>2.5</sub> concentrations across China with  
962 remote sensing, meteorological and land use information, *Sci. Total Environ.*, 636, 52–60,  
963 doi:10.1016/j.scitotenv.2018.04.251, 2018.

964 [Chen, J., Ban, Y., and Li, S.: China: Open access to Earth land-cover map, \*Nature\*, 514\(7523\):](#)  
965 [434-434, doi:10.1038/514434c, 2014.](#)

966 Di, Q., Amini, H., Shi, L., Kloog, I., Silvern, R., Kelly, J., Sabath, M. B., Choirat, C., Koutrakis,  
967 P., Lyapustin, A., Wang, Y., Mickley, L. J. and Schwartz, J.: An ensemble-based model of PM<sub>2.5</sub>  
968 concentration across the contiguous United States with high spatiotemporal resolution, *Environ. Int.*,  
969 130, 104909, doi:10.1016/j.envint.2019.104909, 2019.

970 van Donkelaar, A., Martin, R. V., Brauer, M., Kahn, R., Levy, R., Verduzco, C. and Villeneuve,  
971 P. J.: Global Estimates of Ambient Fine Particulate Matter Concentrations from Satellite-Based  
972 Aerosol Optical Depth: Development and Application, *Environ. Health Perspect.*, 118(6), 847–855,  
973 doi:10.1289/ehp.0901623, 2010.

974 van Donkelaar, A., Martin, R. V., Brauer, M., Hsu, N. C., Kahn, R. A., Levy, R. C., Lyapustin,  
975 A., Sayer, A. M. and Winker, D. M.: Global Estimates of Fine Particulate Matter using a Combined  
976 Geophysical-Statistical Method with Information from Satellites, Models, and Monitors, *Environ. Sci.*  
977 *Technol.*, 50(7), 3762–3772, doi:10.1021/acs.est.5b05833, 2016.

978 Fang, X., Zou, B., Liu, X., Sternberg, T. and Zhai, L.: Satellite-based ground PM<sub>2.5</sub> estimation

979 using timely structure adaptive modeling, *Remote Sens. Environ.*, 186, 152–163,  
980 doi:10.1016/j.rse.2016.08.027, 2016.

981 Fuzzi, S., Baltensperger, U., Carslaw, K., Decesari, S., Denier van der Gon, H., Facchini, M. C.,  
982 Fowler, D., Koren, I., Langford, B., Lohmann, U., Nemitz, E., Pandis, S., Riipinen, I., Rudich, Y.,  
983 Schaap, M., Slowik, J. G., Spracklen, D. V., Vignati, E., Wild, M., Williams, M. and Gilardoni, S.:  
984 Particulate matter, air quality and climate: lessons learned and future needs, *Atmos. Chem. Phys.*,  
985 15(14), 8217–8299, doi:10.5194/acp-15-8217-2015, 2015.

986 Gao, M., Beig, G., Song, S., Zhang, H., Hu, J., Ying, Q., Liang, F., Liu, Y., Wang, H., Lu, X.,  
987 Zhu, T., Carmichael, G. R., Nielsen, C. P. and McElroy, M. B.: The impact of power generation  
988 emissions on ambient PM<sub>2.5</sub> pollution and human health in China and India, *Environ. Int.*,  
989 121(August), 250–259, doi:10.1016/j.envint.2018.09.015, 2018.

990 Geng, G., Xiao, Q., Liu, S., Liu, X., Cheng, J., Zheng, Y., Xue, T., Tong, D., Zheng, B., Peng, Y.,  
991 Huang, X., He, K. and Zhang, Q.: Tracking Air Pollution in China: Near Real-Time PM<sub>2.5</sub> Retrievals  
992 from Multisource Data Fusion, *Environ. Sci. Technol.*, 55(17), 12106–12115,  
993 doi:10.1021/acs.est.1c01863, 2021.

994 Goldberg, D. L., Gupta, P., Wang, K., Jena, C., Zhang, Y., Lu, Z. and Streets, D. G.: Using gap-  
995 filled MAIAC AOD and WRF-Chem to estimate daily PM<sub>2.5</sub> concentrations at 1 km resolution in the  
996 Eastern United States, *Atmos. Environ.*, 199(November 2018), 443–452,  
997 doi:10.1016/j.atmosenv.2018.11.049, 2019.

998 Guo, J., Su, T., Li, Z., Miao, Y., Li, J., Liu, H., Xu, H., Cribb, M. and Zhai, P.: Declining frequency  
999 of summertime local-scale precipitation over eastern China from 1970 to 2010 and its potential link to  
1000 aerosols, *Geophys. Res. Lett.*, 44(11), 5700–5708, doi:10.1002/2017GL073533, 2017.

1001 Hu, X., Waller, L. A., Lyapustin, A., Wang, Y., Al-Hamdan, M. Z., Crosson, W. L., Estes, M. G.,  
1002 Estes, S. M., Quattrochi, D. A., Puttaswamy, S. J. and Liu, Y.: Estimating ground-level PM<sub>2.5</sub>  
1003 concentrations in the Southeastern United States using MAIAC AOD retrievals and a two-stage model,  
1004 *Remote Sens. Environ.*, 140, 220–232, doi:10.1016/j.rse.2013.08.032, 2014.

1005 Huang, C., Hu, J., Xue, T., Xu, H. and Wang, M.: High-Resolution Spatiotemporal Modeling for  
1006 Ambient PM<sub>2.5</sub> Exposure Assessment in China from 2013 to 2019, *Environ. Sci. Technol.*, 55(3),  
1007 2152–2162, doi:10.1021/acs.est.0c05815, 2021.

1008 Kolda, T. G. and Bader, B. W.: Tensor Decompositions and Applications, *SIAM Rev.*, 51(3), 455–

**Deleted:** Folch-Fortuny, A., Arteaga, F. and Ferrer, A.: PCA model building with missing data: New proposals and a comparative study, *Chemom. Intell. Lab. Syst.*, 146, 77–88, doi:10.1016/j.chemolab.2015.05.006, 2015.

**Deleted:** a

**Deleted:** Geng, G., Xiao, Q., Liu, S., Liu, X., Cheng, J., Zheng, Y., Xue, T., Tong, D., Zheng, B., Peng, Y., Huang, X., He, K. and Zhang, Q.: Tracking Air Pollution in China: Near Real-Time PM<sub>2.5</sub> Retrievals from Multisource Data Fusion, *Environ. Sci. Technol.*, acs.est.1c01863, doi:10.1021/acs.est.1c01863, 2021b.

**Deleted:** He, Q., Gu, Y. and Zhang, M.: Spatiotemporal trends of PM<sub>2.5</sub> concentrations in central China from 2003 to 2018 based on MAIAC-derived high-resolution data, *Environ. Int.*, 137(August 2019), 105536, doi:10.1016/j.envint.2020.105536, 2020.

**Deleted:** Jun, C., Ban, Y. and Li, S.: Open access to Earth land-cover map, *Nature*, 514(7523), 434–434, doi:10.1038/514434c, 2014.

1028 500, doi:10.1137/07070111X, 2009.

1029 de Leeuw, G., Sogacheva, L., Rodriguez, E., Kourtidis, K., Georgoulas, A. K., Alexandri, G.,  
1030 Amiridis, V., Proestakis, E., Marinou, E., Xue, Y. and van der A, R.: Two decades of satellite  
1031 observations of AOD over mainland China using ATSR-2, AATSR and MODIS/Terra: data set  
1032 evaluation and large-scale patterns, *Atmos. Chem. Phys.*, 18(3), 1573–1592, doi:10.5194/acp-18-  
1033 1573-2018, 2018.

1034 Li, J., Li, C. and Zhao, C.: Different trends in extreme and median surface aerosol extinction  
1035 coefficients over China inferred from quality-controlled visibility data, *Atmos. Chem. Phys.*, 18(5),  
1036 3289–3298, doi:10.5194/acp-18-3289-2018, 2018a.

1037 ~~Li, L., Franklin, M., Girguis, M., Lurmann, F., Wu, J., Pavlovic, N., Breton, C., Gilliland, F. and~~  
1038 ~~Habre, R.: Spatiotemporal imputation of MAIAC AOD using deep learning with downscaling, *Remote*~~  
1039 ~~*Sens. Environ.*, 237(October 2019), 111584, doi:10.1016/j.rse.2019.111584, 2020.~~

1040 Li, K., Bai, K., Li, Z., Guo, J. and Chang, N.-B.: Synergistic Data Fusion of Multimodal AOD and  
1041 Air Quality Data for Near Real-Time Full Coverage Air Pollution Assessment, *J. Environ. Manage.*,  
1042 ~~302, 114121, doi: 10.1016/j.jenvman.2021.114121, 2022.~~

1043 Li, Z., Zhang, Y., Shao, J., Li, B., Hong, J., Liu, D., Li, D., Wei, P., Li, W., Li, L., Zhang, F., Guo,  
1044 J., Deng, Q., Wang, B., Cui, C., Zhang, W., Wang, Z., Lv, Y., Xu, H., Chen, X., Li, L. and Qie, L.:  
1045 Remote sensing of atmospheric particulate mass of dry PM<sub>2.5</sub> near the ground: Method validation  
1046 using ground-based measurements, *Remote Sens. Environ.*, 173, 59–68, doi:10.1016/j.rse.2015.11.019,  
1047 2016.

1048 Li, Z., Wang, Y., Guo, J., Zhao, C., Cribb, M. C., Dong, X., Fan, J., Gong, D., Huang, J., Jiang,  
1049 M., Jiang, Y., Lee, S. S., Li, H., Li, J., Liu, J., Qian, Y., Rosenfeld, D., Shan, S., Sun, Y., Wang, H.,  
1050 Xin, J., Yan, X., Yang, X., Yang, X. qun, Zhang, F. and Zheng, Y.: East Asian Study of Tropospheric  
1051 Aerosols and their Impact on Regional Clouds, Precipitation, and Climate (EAST-AIRCPC), *J.*  
1052 *Geophys. Res. Atmos.*, 124(23), 13026–13054, doi:10.1029/2019JD030758, 2019.

1053 Lin, C., Li, Y., Lau, A. K. H., Deng, X., Tse, T. K. T., Fung, J. C. H., Li, C., Li, Z., Lu, X., Zhang,  
1054 X. and Yu, Q.: Estimation of long-term population exposure to PM<sub>2.5</sub> for dense urban areas using 1-  
1055 km MODIS data, *Remote Sens. Environ.*, 179, 13–22, doi:10.1016/j.rse.2016.03.023, 2016.

1056 Liu, M., Bi, J. and Ma, Z.: Visibility-Based PM<sub>2.5</sub> Concentrations in China: 1957–1964 and  
1057 1973–2014, *Environ. Sci. Technol.*, 51(22), 13161–13169, doi:10.1021/acs.est.7b03468, 2017.

**Deleted:** Li, L., Zhang, J., Meng, X., Fang, Y., Ge, Y., Wang, J., Wang, C., Wu, J. and Kan, H.: Estimation of PM<sub>2.5</sub> concentrations at a high spatiotemporal resolution using constrained mixed-effect bagging models with MAIAC aerosol optical depth, *Remote Sens. Environ.*, 217(January), 573–586, doi:10.1016/j.rse.2018.09.001, 2018b.

**Deleted:** 2021. In revision

1065 Liu, Y., Paciorek, C. J. and Koutrakis, P.: Estimating Regional Spatial and Temporal Variability  
1066 of PM<sub>2.5</sub> Concentrations Using Satellite Data, Meteorology, and Land Use Information, Environ.  
1067 Health Perspect., 117(6), 886–892, doi:10.1289/ehp.0800123, 2009.

1068 Lyapustin, A., Martonchik, J., Wang, Y., Laszlo, I. and Korin, S.: Multiangle implementation of  
1069 atmospheric correction (MAIAC): 1. Radiative transfer basis and look-up tables, J. Geophys. Res.  
1070 Atmos., 116(3), doi:10.1029/2010JD014985, 2011.

1071 Lyapustin, A., Wang, Y., Korin, S. and Huang, D.: MODIS Collection 6 MAIAC algorithm,  
1072 Atmos. Meas. Tech., 11(10), 5741–5765, doi:10.5194/amt-11-5741-2018, 2018.

1073 Lyu, B., Hu, Y., Zhang, W., Du, Y., Luo, B., Sun, X., Sun, Z., Deng, Z., Wang, X., Liu, J., Wang,  
1074 X. and Russell, A. G.: Fusion Method Combining Ground-Level Observations with Chemical  
1075 Transport Model Predictions Using an Ensemble Deep Learning Framework: Application in China to  
1076 Estimate Spatiotemporally-Resolved PM<sub>2.5</sub> Exposure Fields in 2014–2017, Environ. Sci. Technol.,  
1077 53(13), 7306–7315, doi:10.1021/acs.est.9b01117, 2019.

1078 Ma, Z., Hu, X., Sayer, A. M., Levy, R., Zhang, Q., Xue, Y., Tong, S., Bi, J., Huang, L. and Liu,  
1079 Y.: Satellite-based spatiotemporal trends in PM<sub>2.5</sub> concentrations: China, 2004–2013, Environ. Health  
1080 Perspect., 124(2), 184–192, doi:10.1289/ehp.1409481, 2016.

1081 Park, S., Lee, J., Im, J., Song, C. K., Choi, M., Kim, J., Lee, S., Park, R., Kim, S. M., Yoon, J.,  
1082 Lee, D. W. and Quackenbush, L. J.: Estimation of spatially continuous daytime particulate matter  
1083 concentrations under all sky conditions through the synergistic use of satellite-based AOD and  
1084 numerical models, Sci. Total Environ., 713, 136516, doi:10.1016/j.scitotenv.2020.136516, 2020.

1085 Shen, F., Zhang, L., Jiang, L., Tang, M., Gai, X., Chen, M. and Ge, X.: Temporal variations of six  
1086 ambient criteria air pollutants from 2015 to 2018, their spatial distributions, health risks and  
1087 relationships with socioeconomic factors during 2018 in China, Environ. Int., 137(February), 105556,  
1088 doi:10.1016/j.envint.2020.105556, 2020.

1089 Sidiropoulos, N. D., De Lathauwer, L., Fu, X., Huang, K., Papalexakis, E. E. and Faloutsos, C.:  
1090 Tensor Decomposition for Signal Processing and Machine Learning, IEEE Trans. Signal Process.,  
1091 65(13), 3551–3582, doi:10.1109/TSP.2017.2690524, 2017.

1092 Sogacheva, L., Popp, T., Sayer, A. M., Dubovik, O., Garay, M. J., Heckel, A., Christina Hsu, N.,  
1093 Jethva, H., Kahn, R. A., Kolmonen, P., Kosmale, M., De Leeuw, G., Levy, R. C., Litvinov, P.,  
1094 Lyapustin, A., North, P., Torres, O. and Arola, A.: Merging regional and global aerosol optical depth

**Deleted:** Ma, Z., Hu, X., Huang, L., Bi, J. and Liu, Y.:  
Estimating Ground-Level PM<sub>2.5</sub> in China Using Satellite  
Remote Sensing, Environ. Sci. Technol., 48(13), 7436–7444,  
doi:10.1021/es5009399, 2014.

1099 records from major available satellite products, *Atmos. Chem. Phys.*, 20(4), 2031–2056,  
1100 doi:10.5194/acp-20-2031-2020, 2020.

1101 Sun, J.-L., Jing, X., Chang, W.-J., Chen, Z.-X. and Zeng, H.: Cumulative health risk assessment  
1102 of halogenated and parent polycyclic aromatic hydrocarbons associated with particulate matters in  
1103 urban air, *Ecotoxicol. Environ. Saf.*, 113, 31–37, doi:10.1016/j.ecoenv.2014.11.024, 2015.

1104 Sun, Z., Chang, N. Bin, Chen, C. F., Mostafiz, C. and Gao, W.: Ensemble learning via higher  
1105 order singular value decomposition for integrating data and classifier fusion in water quality  
1106 monitoring, *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, 14, 3345–3360,  
1107 doi:10.1109/JSTARS.2021.3055798, 2021.

1108 Tang, Q., Bo, Y. and Zhu, Y.: Spatiotemporal fusion of multiple-satellite aerosol optical depth  
1109 (AOD) products using Bayesian maximum entropy method, *J. Geophys. Res. Atmos.*, 121(8), 4034–  
1110 4048, doi:10.1002/2015JD024571, 2016.

1111 Tucker, L. R.: Some mathematical notes on three-mode factor analysis, *Psychometrika*, 31(3),  
1112 279–311, doi:10.1007/BF02289464, 1966.

1113 Wang, B., Yuan, Q., Yang, Q., Zhu, L., Li, T. and Zhang, L.: Estimate hourly PM<sub>2.5</sub>  
1114 concentrations from Himawari-8 TOA reflectance directly using geo-intelligent long short-term  
1115 memory network, *Environ. Pollut.*, 271, 116327, doi:10.1016/j.envpol.2020.116327, 2021a.

1116 Wang, Q., Shen, Y., and Zhang, J. Q.: A nonlinear correlation measure for multivariable data  
1117 set, *Phys. D*, 3–4, 287–295, doi:10.1016/j.physd.2004.11.001, 2005.

1118 Wang, Y., Yuan, Q., Li, T., Shen, H., Zheng, L. and Zhang, L.: Large-scale MODIS AOD products  
1119 recovery: Spatial-temporal hybrid fusion considering aerosol variation mitigation, *ISPRS J.*  
1120 *Photogramm. Remote Sens.*, 157(July), 1–12, doi:10.1016/j.isprsjprs.2019.08.017, 2019.

1121 Wang, Y., Yuan, Q., Li, T., Tan, S. and Zhang, L.: Full-coverage spatiotemporal mapping of  
1122 ambient PM<sub>2.5</sub> and PM<sub>10</sub> over China from Sentinel-5P and assimilated datasets: Considering the  
1123 precursors and chemical compositions, *Sci. Total Environ.*, 793, 148535,  
1124 doi:10.1016/j.scitotenv.2021.148535, 2021b.

1125 Wei, J., Li, Z., Peng, Y. and Sun, L.: MODIS Collection 6.1 aerosol optical depth products over  
1126 land and ocean: validation and comparison, *Atmos. Environ.*, 201, 428–440,  
1127 doi:10.1016/j.atmosenv.2018.12.004, 2019b.

1128 Wei, J., Li, Z., Lyapustin, A., Sun, L., Peng, Y., Xue, W., Su, T. and Cribb, M.: Reconstructing 1-

**Deleted:** Wei, J., Huang, W., Li, Z., Xue, W., Peng, Y., Sun, L. and Cribb, M.: Estimating 1-km-resolution PM<sub>2.5</sub> concentrations across China using the space-time random forest approach, *Remote Sens. Environ.*, 231(May), 111221, doi:10.1016/j.rse.2019.111221, 2019a.¶

**Deleted:** Wei, J., Li, Z., Cribb, M., Huang, W., Xue, W., Sun, L., Guo, J., Peng, Y., Li, J., Lyapustin, A., Liu, L., Wu, H. and Song, Y.: Improved 1 km resolution PM<sub>2.5</sub> estimates across China using enhanced space – time extremely randomized trees, *Atmos. Chem. Phys.*, 20, 3273–3289, 2020a.¶

1140 km-resolution high-quality PM2.5 data records from 2000 to 2018 in China: spatiotemporal variations  
1141 and policy implications, *Remote Sens. Environ.*, 252(January 2020), 112136,  
1142 doi:10.1016/j.rse.2020.112136, 2021a.

1143 Wei, X., Chang, N., Bai, K. and Gao, W.: Satellite remote sensing of aerosol optical depth:  
1144 advances, challenges, and perspectives, *Crit. Rev. Environ. Sci. Technol.*, 50(16), 1640–1725,  
1145 doi:10.1080/10643389.2019.1665944, 2020.

1146 Wei, X., Bai, K., Chang, N. and Gao, W.: Multi-source hierarchical data fusion for high-resolution  
1147 AOD mapping in a forest fire event, *Int. J. Appl. Earth Obs. Geoinf.*, 102(May), 102366,  
1148 doi:10.1016/j.jag.2021.102366, 2021b.

1149 Xiao, Q., Zhang, H., Choi, M., Li, S., Kondragunta, S., Kim, J., Holben, B., Levy, R. C. and Liu,  
1150 Y.: Evaluation of VIIRS, GOCI, and MODIS Collection 6 AOD retrievals against ground  
1151 sunphotometer observations over East Asia, *Atmos. Chem. Phys.*, 16(3), 1255–1269, doi:10.5194/acp-  
1152 16-1255-2016, 2016.

1153 Xiao, Q., Wang, Y., Chang, H. H., Meng, X., Geng, G., Lyapustin, A. and Liu, Y.: Full-coverage  
1154 high-resolution daily PM2.5 estimation using MAIAC AOD in the Yangtze River Delta of China,  
1155 *Remote Sens. Environ.*, 199(May), 437–446, doi:10.1016/j.rse.2017.07.023, 2017.

1156 Xiao, Q., Chang, H. H., Geng, G. and Liu, Y.: An Ensemble Machine-Learning Model to Predict  
1157 Historical PM2.5 Concentrations in China from Satellite Data, *Environ. Sci. Technol.*,  
1158 doi:10.1021/acs.est.8b02917, 2018.

1159 Xin, J., Wang, Y., Pan, Y., Ji, D., Liu, Z., Wen, T., Wang, Y., Li, X., Sun, Y., Sun, J., Wang, P.,  
1160 Wang, G., Wang, X., Cong, Z., Song, T., Hu, B., Wang, L., Tang, G., Gao, W., Guo, Y., Miao, H.,  
1161 Tian, S. and Wang, L.: The Campaign on Atmospheric Aerosol Research Network of China: CARE-  
1162 China, *Bull. Am. Meteorol. Soc.*, 96(7), 1137–1155, doi:10.1175/BAMS-D-14-00039.1, 2015.

1163 Xu, H., Guang, J., Xue, Y., de Leeuw, G., Che, Y. H., Guo, J., He, X. W. and Wang, T. K.: A  
1164 consistent aerosol optical depth (AOD) dataset over mainland China by integration of several AOD  
1165 products, *Atmos. Environ.*, 114, 48–56, doi:10.1016/j.atmosenv.2015.05.023, 2015.

1166 Xue, T., Zheng, Y., Tong, D., Zheng, B., Li, X., Zhu, T. and Zhang, Q.: Spatiotemporal continuous  
1167 estimates of PM2.5 concentrations in China, 2000–2016: A machine learning method with inputs from  
1168 satellites, chemical transport model, and ground observations, *Environ. Int.*, 123(December 2018),  
1169 345–357, doi:10.1016/j.envint.2018.11.075, 2019.

Deleted: b

Deleted: a

Deleted: Xiao, Q., Wang, Y., Chang, H. H., Meng, X., Geng, G., Lyapustin, A. and Liu, Y.: Full-coverage high-resolution daily PM2.5 estimation using MAIAC AOD in the Yangtze River Delta of China, *Remote Sens. Environ.*, 199, 437–446, doi:10.1016/j.rse.2017.07.023, 2017b.

Deleted: Xiao, Q., Geng, G., Liang, F., Wang, X., Lv, Z., Lei, Y., Huang, X., Zhang, Q., Liu, Y. and He, K.: Changes in spatial patterns of PM2.5 pollution in China 2000–2018: Impact of clean air policies, *Environ. Int.*, 141(April), 105776, doi:10.1016/j.envint.2020.105776, 2020.

1182 Yang, F., Tan, J., Zhao, Q., Du, Z., He, K., Ma, Y., Duan, F., Chen, G. and Zhao, Q.:  
1183 Characteristics of PM<sub>2.5</sub> speciation in representative megacities and across China, *Atmos. Chem.*  
1184 *Phys.*, 11(11), 5207–5219, doi:10.5194/acp-11-5207-2011, 2011.

1185 Yang, J. and Huang, X.: 30 m annual land cover and its dynamics in China from 1990 to 2019,  
1186 *Earth Syst. Sci. Data*, 13, 3907–3925, doi: 10.5194/essd-13-3907, 2021.

1187 [Yang Y., Zheng Z., Yim S.H.L., Roth M., Ren G., Gao Z., Wang T., Li Q., Shi C., Ning G. and](#)  
1188 [Li Y.B.: PM<sub>2.5</sub> Pollution Modulates Wintertime Urban-Heat-Island Intensity in the Beijing-Tianjin-](#)  
1189 [Hebei Megalopolis, China. \*Geophys. Res. Lett.\*, 47\(1\), e2019GL084288, doi:10.1029/2019gl084288,](#)  
1190 [2020.](#)

1191 [Zhang, T., Zeng, C., Gong, W., Wang, L., Sun, K., Shen, H., Zhu, Z. and Zhu, Z.:](#) Improving  
1192 spatial coverage for Aqua MODIS AOD using NDVI-based multi-temporal regression analysis,  
1193 *Remote Sens.*, 9(4), doi:10.3390/rs9040340, 2017.

1194 Zhang, Y., Gao, L., Cao, L., Yan, Z. and Wu, Y.: Decreasing atmospheric visibility associated  
1195 with weakening winds from 1980 to 2017 over China, *Atmos. Environ.*, 224(July 2019), 117314,  
1196 doi:10.1016/j.atmosenv.2020.117314, 2020.

1197 Zhao, C., Yang, Y., Fan, H., Huang, J., Fu, Y., Zhang, X., Kang, S., Cong, Z., Letu, H. and Menenti,  
1198 M.: Aerosol characteristics and impacts on weather and climate over the Tibetan Plateau, *Natl. Sci.*  
1199 *Rev.*, 7(3), 492–495, doi:10.1093/nsr/nwz184, 2020.

1200 [Zheng, Z., Ren, G., Wang, H., Dou, J., Gao, Z., Duan, C., Li, Y. and Ngarukiyimana, J., Zhao,](#)  
1201 [C., Cao, C., Jiang, M., and Yang, Y.: Relationship between Fine Particle Pollution and the Urban Heat](#)  
1202 [Island in Beijing, China: Observational Evidence, \*Boundary-Layer Meteorol.\*, 169\(1\), 93-113, doi:](#)  
1203 [10.1007/s10546-018-0362-6, 2018.](#)

1204 [Zheng Z., Zhao, C., Lolli, S., Wang, X., Wang, Y., Ma, X., Li, Q. and Yang, Y.:](#) Diurnal Variation  
1205 of Summer Precipitation Modulated by Air Pollution: Observational Evidences in the Beijing  
1206 Metropolitan Area, *Environ. Res. Lett.*, 15, 094053, doi:10.1088/1748-9326/ab99fc, 2020.

1207  
1208

Deleted: Discuss., 2021(April), 1–29,  
doi:https://doi.org/10.5194/essd-2021-7, 2021.

Formatted: Subscript

Deleted: Yi-Lei Chen, Chiou-Ting Hsu and Liao, H.-Y. M.:  
Simultaneous Tensor Decomposition and Completion Using  
Factor Priors, *IEEE Trans. Pattern Anal. Mach. Intell.*, 36(3),  
577–591, doi:10.1109/TPAMI.2013.164, 2014.

Formatted: Indent: First line: 0 cm