





18

### Abstract

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



## 48 1 Introduction

49 Atmospheric aerosols not only impact regional climate by changing the Earth radiation budget  
50 but significantly influence air quality at the ground level (Fuzzi et al., 2015; Gao et al., 2018; Shen et  
51 al., 2020; Sun et al., 2015). Monitoring aerosol loading in the atmosphere is thus of great significance  
52 for climate change attribution and haze pollution assessment. Aerosol optical depth (AOD), a measure  
53 of aerosols distributed within a column of air from the Earth's surface to the top of the atmosphere,  
54 has been monitored for decades to quantify aerosol loading in the atmosphere. Compared with sparsely  
55 distributed ground aerosol monitoring stations (e.g., AERONET), satellite instruments can provide  
56 better AOD observations because of vast spatial coverage and high sampling frequency. An overview  
57 of sensors, algorithms, and AOD datasets that are widely used can be found in the literature such as  
58 Sogacheva et al. (2020) and Wei et al. (2020).

59 Due to negative impacts of bright surface (e.g., snow cover) and clouds, as well as algorithmic  
60 restrictions, satellite AOD retrievals often suffer from extensive data gaps, significantly reducing the  
61 downstream application potential such as mapping particulate matter (PM) concentrations at the  
62 ground surface (e.g., Bai et al., 2019; Wei et al., 2021a). Also, data gaps in AOD imageries may result  
63 in large uncertainty when assessing aerosol impacts on weather and climate (Guo et al., 2017; Li et al.,  
64 2019; Zhao et al., 2020). Over the years, many gap filling methods have been developed (e.g., Bai et  
65 al., 2016, 2020b; Chang et al., 2015). Nonetheless, filling data gaps in satellite-based AOD products  
66 is still a challenge due to extraordinary nonrandom missing values and high aerosol dynamics in space  
67 and time.

68 Wei et al. (2020a) provided a short review of methods that have been frequently applied to deal  
69 with data gaps in AOD products. In general, merging AOD data acquired from diverse instruments  
70 and/or platforms is the most popular approach to improve AOD spatial coverage (Sogacheva et al.,  
71 2020). Statistical methods such as linear regression (Bai et al., 2019a; Wang et al., 2019; Zhang et al.,  
72 2017), inversed variance weighting (Chen et al., 2018; Ma et al., 2016; Sogacheva et al., 2020), and  
73 maximum likelihood estimate (Xu et al., 2015), are often applied to account for systematic bias among  
74 different datasets. Data fusion methods such as Bayesian maximum entropy could be applied to blend  
75 AOD products with different resolutions (Tang et al., 2016; Wei et al., 2021b). Another way is to  
76 reconstruct missing AOD values using either neighboring observations in space and time or external



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

79         Although there exist many versatile gap filling methods, spatially gap free AOD datasets are  
80 always rare, particularly satellite-based high-resolution AOD datasets, resulting in significant limit in  
81 downstream applications such as  $PM_x$  concentration mapping. In spite of versatile  $PM_{2.5}$  concentration  
82 prediction models (e.g., Di et al., 2019; Fang et al., 2016; He et al., 2020; Hu et al., 2014; Li et al.,  
83 2018b, 2016; Lin et al., 2016; Liu et al., 2009; Ma et al., 2014; Wang et al., 2021a), to date, there are  
84 few publicly accessible PM concentration datasets that can be used to examine haze pollution  
85 variations regionally and globally. As typical datasets, the one generated by the Dalhousie University  
86 (van Donkelaar et al., 2010, 2016), CHAP (Wei et al., 2019a), and TAP (Geng et al., 2021a)  
87 demonstrated the global effort to elevate earth system science research. However, these datasets more  
88 or less suffer from drawbacks in spatial and/or temporal resolution, spatial coverage, and data accuracy.  
89 To meet the contemporary needs, Zhang et al. (2021) provided a more comprehensive review of the  
90 widely used PM concentration mapping approaches.

91         With a thorough review for  $PM_{2.5}$  concentration mapping techniques, an optimal full-coverage  
92  $PM_{2.5}$  concentration mapping scheme was proposed, in which diverse aerosol datasets were fused  
93 toward a full-coverage AOD map based on a multi-modal approach (Bai et al., 2021). In parallel with  
94 these efforts, some attempted to improve AOD data coverage over space with high accuracy by  
95 merging AODs observed at adjacent times directly (Li et al., 2021). Given such prior knowledge, the  
96 current study developed a big data analytics framework for generating a Long-term Gap-free High-  
97 resolution Air Pollutants concentration dataset (abbreviated as LGHAP hereafter) providing AOD,  
98  $PM_{2.5}$  and  $PM_{10}$  concentration with a daily 1-km resolution in China from 2000 to 2020. Multimodal  
99 aerosol data acquired from diverse sources including satellites, ground stations and numerical models  
100 were synergistically integrated via the higher order singular value decomposition (HOSVD) to form a  
101 tensor flow based data fusion method in the current study. Full coverage  $PM_{2.5}$  and  $PM_{10}$  concentration  
102 data were then estimated on the basis of the gap-filled AOD dataset. This 21-year-long gap-free high  
103 resolution (daily/1km) aerosol dataset was then compared against ground-based AOD and PM  
104 observations to evaluate the data accuracy of each product, particularly the performance in spatial  
105 pattern recognition and temporal trend assessment. These advances led to explore the long-term



106 variability and population exposure to haze pollution in China over the past two decades by taking  
 107 advantage of the LGHAP dataset.

108 **2 Data sources**

109 Table 1 summarizes the multisource datasets used in this study to help generate the LGHAP  
 110 product. As shown, six satellite-based AOD products, five numerical simulations of AOD and aerosol  
 111 components, eleven meteorological factors, six ground-based AOD and air quality datasets, as well as  
 112 five land cover, topographic and socioeconomic parameters, were employed. Descriptions of these  
 113 datasets are given in the following subsections.

114 **Table 1.** Summary of the data sources used in this study to generate gap free high resolution AOD  
 115 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°
	AERONET	2000–2020	hourly	point
Meteorology	Air temperature		hourly	0.25°
	U/V component of wind		hourly	0.25°
	Relative humidity		hourly	0.25°
	Surface pressure	2000–2020	hourly	0.25°
	Boundary layer height		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°



## 116 2.1 Gridded aerosol products

117 In many previous studies, coarse AOD and/or aerosol components simulations acquired from  
118 numerical models were oftentimes used as the primary data source to help derive full-coverage AOD  
119 and/or PM<sub>2.5</sub> concentration maps (e.g., Park et al., 2020; Wang et al., 2021b). However, due to the lack  
120 of high accuracy near real-time emission inventory, simulated AOD and/or aerosol components are  
121 often prone to large uncertainty, which could be inevitably introduced to the final PM<sub>2.5</sub> estimations if  
122 no observational data are applied for bias correction. In such a research context, here we used six  
123 satellite-based AOD products with a relatively long-term coverage to help better reconstruct historical  
124 AOD variations over space and time.

125 The latest AOD product derived from the MODerate-resolution Imaging Spectroradiometer  
126 (MODIS) onboard Terra using the multiangle implementation of atmospheric correction (MAIAC)  
127 algorithm (Lyapustin et al., 2011, 2018), was hereby used as the baseline dataset for the generation of  
128 full-coverage AOD maps. This AOD product has not only a finer spatial resolution but a comparable  
129 and even better accuracy, when comparing with those derived from the Dark Target and Deep Blue  
130 algorithms (Goldberg et al., 2019; Lyapustin et al., 2018; Xiao et al., 2017b). In addition, AOD  
131 products derived from MODIS onboard Aqua, the Multi-angle Imaging SpectroRadiometer (MISR)  
132 onboard Terra, Visible Infrared Imaging Radiometer Suite (VIIRS) onboard Suomi-NPP, Advanced  
133 Along-Track Scanning Radiometer (AATSR) onboard Envisat and POLarization and Directionality of  
134 the Earth's Reflectances (POLDER) onboard PARASOL, were also utilized. The ultimate goal was to  
135 reduce bias in full-coverage AOD imagery by better spatial coverage of observational AOD as much  
136 as possible. Accuracies of these AOD products have been extensively validated, e.g., de Leeuw et al.  
137 (2018), Xiao et al. (2016), Wei et al. (2019b), Che et al. (2019), to name a few, in the current study. A  
138 brief description of these satellite-based AOD products can be found in Text S1 in the supplementary  
139 information.

140 In addition to satellite-based AOD products, numerically simulated aerosol diagnostics from  
141 MERRA-2, including AOD and aerosol components such as black carbon, organic carbon, dust and  
142 sulfate, were also applied to help reconstruct missing AOD information and to predict PM<sub>2.5</sub> and PM<sub>10</sub>  
143 concentrations at the ground level. The aerosol components were used here as a proxy of emission  
144 inventory when predicting PM concentrations. Big data analytics procedures applied to these datasets  
145 will be described in section 3.



## 146 2.2 *In situ* AOD and air quality measurements

147 AOD observations from Aerosol Robotic Network (AERONET) were used to evaluate the  
148 prediction accuracy of the generated full-coverage (gap free) AOD product, as well as the learning  
149 target to infer AOD from air pollutants concentration and atmospheric visibility. Considering few valid  
150 data were provided in the Level 2.0 dataset, here we used the Level 1.5 AOD data to guarantee adequate  
151 *in situ* AOD data coverage in space and time. To validate the gridded AOD products in this study, each  
152 *in situ* AOD observation was registered with the gridded mean AOD over a 50×50 km window.

153 Near-surface air pollutants concentrations including PM<sub>2.5</sub>, PM<sub>10</sub>, NO<sub>2</sub>, and SO<sub>2</sub> that were  
154 sampled at state-controlled monitoring sites were also applied, not only to help establish machine-  
155 learned regression models for PM prediction (PM<sub>2.5</sub> and PM<sub>10</sub>), but to infer AOD over air quality  
156 monitoring sites given their dense distribution across China. The gauged air pollutants concentration  
157 data have been released online on an hourly basis by the China National Environment Monitoring  
158 Center since the late 2013. For quality control, outliers were first detected and removed from each  
159 pollutant dataset by following the criteria used in our previous study (Bai et al., 2020a). The missing  
160 values were then reconstructed using the diurnal cycle constrained empirical orthogonal function  
161 (DCCEOF) method proposed in Bai et al. (2020b).

162 The 3-hour resolution atmospheric visibility data acquired from 4,052 weather stations were  
163 employed to help generate gap free AOD maps before 2014, at which *in situ* air quality measurements  
164 were not available. Previous studies have attempted to predict PM<sub>2.5</sub> concentration from atmospheric  
165 visibility data with good accuracies (Liu et al., 2017), indicative of a great potential for estimating  
166 AOD. Specifically, visibility data were used as an important predictor for site-specific AOD prediction,  
167 and the resulting AOD predictions were then used as a critical prior information for reconstructing  
168 AOD distributions over space, especially over those regions without satellite AOD observations. Since  
169 automatic visibility sensors have been widely used at many sites since 2014, those data were excluded  
170 to guarantee the data consistency (Li et al., 2018a). For quality control, the consistency of visibility  
171 data was examined using an outlier detection method, i.e., the annual mean should not exceed 3 times  
172 the standard deviation of data over a 5-year time window (Zhang et al., 2020). Those with apparent  
173 jumps and drifts in visibility time series were excluded. Meanwhile, visibility data on rainstorm and  
174 foggy days were eliminated as well.



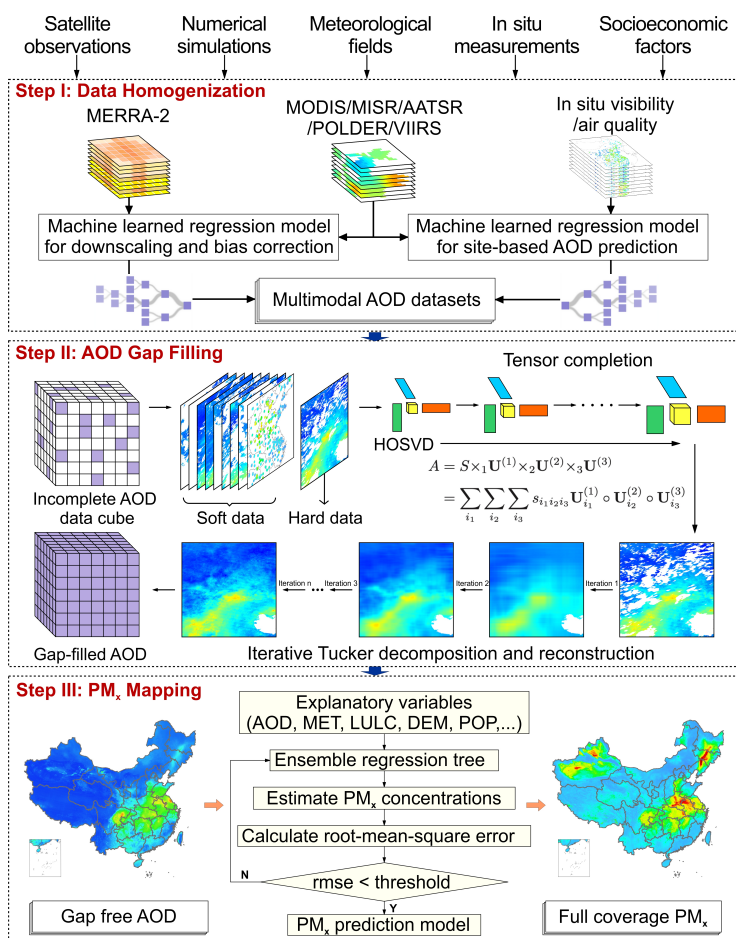
### 175 **2.3 Auxiliary data**

176 As shown in Table 1, eleven meteorological factors, including air temperature at the near surface,  
177 wind speed and direction, relative humidity, surface pressure, boundary layer height, total column  
178 water vapor, downwards solar radiation, and instantaneous moisture flux, were used to help infer PM<sub>2.5</sub>  
179 and PM<sub>10</sub> from AOD, as well as to downscale AOD from MERRA-2. These data were acquired from  
180 the fifth generation ECMWF atmospheric reanalysis (ERA-5), and the first three factors were extracted  
181 at the levels of not only ground surface but 850 hpa and 500 hpa so as to indicate the vertical structure  
182 of the atmosphere. Additionally, population data from WorldPop, land cover from CLCD during 2000  
183 to 2019 (Yang and Huang, 2021) and GLOBELAND 30 in 2020 (Chen et al., 2014), elevation data  
184 from the Global Digital Elevation Model (GDEM) version 2, as well as monthly composited 1-km  
185 normalized difference vegetation index (NDVI) from MODIS, were employed to indicate the  
186 socioeconomic and ecological contributions to haze pollutions. Properties of these datasets can be  
187 found in Table 1, and datasets with a finer resolution were upscaled to 0.01° via a cubic interpolation  
188 method.

### 189 **3 Methodology**

190 To advance environment management and earth system science analysis, the current study  
191 developed a big data analytics framework for generating long-term gap free aerosol spatiotemporal  
192 datasets and demonstrated its applications in China. Such big data analytics was constructed via a  
193 seamless integration of the tensor flow based multimodal data fusion with ensemble learning based  
194 PM concentration estimation. When generating this dataset, the proposed method transformed a set of  
195 data tensors of AOD and other related datasets such as air pollutants concentration and atmospheric  
196 visibility that were acquired from diversified sensors or platforms via integrative efforts of spatial  
197 pattern recognition for high dimensional gridded data analysis toward data fusion and multiresolution  
198 image analysis, as well as knowledge transfer in statistical data mining. The proposed method consists  
199 of three major procedures in general, including multisensory data homogenization, tensor flow based  
200 AOD reconstruction, and ensemble learning for PM concentration estimation. The analytical  
201 framework of the big data analytics is depicted in Figure 1 and described in details in the following  
202 subsections.





203

204 **Figure 1.** Flowchart of the proposed big data analytics framework for generating a long-term gap-free  
 205 high-resolution air pollutants concentration dataset (LGHAP), taking AOD and PM concentration in  
 206 China as illustration.

### 207 3.1 Multisensory data homogenization

208 Since a set of aerosol products with different types, resolution, and accuracies were applied to  
 209 support the generation of gap-free AOD imageries, harmonizing cross-mission biases and scale  
 210 differences between these diversified datasets is thus of critical importance to facilitate multisensory  
 211 data integration. In this study, machine-learned regression models were established to harmonize these  
 212 heterogeneous aerosol datasets. A baseline dataset was first selected to be used as the learning target  
 213 while other datasets were calibrated to the level of baseline dataset to make them comparable. Given



214 finer resolution and higher proportion of data coverage in space and time, the MAIAC AOD product  
215 from Terra ( $AOD_{Terra}$ ) was selected as the baseline dataset. Consequently, six machine-learned  
216 regression models were established between  $AOD_{Terra}$  and each gridded AOD product (i.e., five  
217 satellite-based AOD products plus MERRA-2 AOD simulations) using the random forest method.  
218 Meteorological factors, land cover, topographic and socioeconomic variables were used as covariates  
219 to help downscale these multimodal AOD products to have a resolution same as  $AOD_{Terra}$  while  
220 accounting for cross-mission biases arising from temporal and algorithmic differences.

221 Considering data gaps are extensive in satellite-based AOD products, especially over regions  
222 with thick cloud cover, providing prior AOD information over such region is thus of great value in  
223 support of the reconstruction of missing AOD values. As indicated in our recent studies, AOD can be  
224 accurately predicted from ground measured air pollutants concentration, showing an accuracy even  
225 over some satellite AOD retrievals (Li et al., 2021; Bai et al., 2021). To support AOD reconstruction  
226 over regions without satellite AOD observations, we attempted to infer AOD over air quality  
227 monitoring sites from air pollutants concentration measurements via an ensemble learning approach.  
228 Similarly, machine-learned regression models were established using random forest by taking  
229  $AOD_{Terra}$  as the learning target while ground measured air pollutants concentration, meteorological  
230 factors, land cover, and terrain information, were used as predictors.

231 The transformation of ground measured air pollutants concentration data to AOD empowers us  
232 to provide external observational AOD to supplement satellite observations, especially over regions  
233 suffering from significant data gaps. Since air pollutants concentration data were not available before  
234 2013, atmospheric visibility data sampled at dense weather stations were hereby used as an alternative  
235 for AOD prediction, by applying a similar prediction model as used above for air pollutants  
236 concentration. Figure S1 show the ground-based validation results of AOD inferred from atmospheric  
237 visibility and air pollutants concentration, indicative of a generally good accuracy of these inferred  
238 AOD values. All efforts led to aggregate a set of multimodal aerosol data with different properties for  
239 multisensory data fusion toward generating full-coverage (gap free) AOD mapping as the next step.

### 240 **3.2 Tensor flow based AOD reconstruction**

241 The core of generating full coverage AOD imageries is to fill in data gaps in  $AOD_{Terra}$ . Previous  
242 studies have demonstrated that merging satellite AOD retrievals derived at adjacent time steps can



243 help improve the observational AOD coverage at each single snapshot, while the involvement of  
244 numerical AOD simulations can help fill in AOD data gaps (Li et al., 2021; Bai et al., 2021). In this  
245 study, a tensor completion method was particularly designed and applied to fulfil the gap filling in  
246 AOD<sub>Terra</sub>. Specifically, the incomplete AOD<sub>Terra</sub> imageries were deemed as the hard data (true AOD  
247 state) while other AODs datasets (e.g., the downscaled AOD datasets and site-specific AOD  
248 predictions inferred from air pollutants concentration and atmospheric visibility) were used as the soft  
249 data to help reconstruct AOD distribution in AOD<sub>Terra</sub> via tensor flow based pattern recognition.  
250 Detailed procedures for gap filling are outlined as follows.

### 251 3.2.1 Initial AOD tensor construction

252 Due to extensive data gaps in satellite-based AOD retrievals, it is insufficient to reconstruct all  
253 missing AOD information in AOD<sub>Terra</sub> for a given date by simply using the harmonized satellite-based  
254 AOD data synchronously. To fulfill AOD gap filling, the newly developed tensor completion method  
255 was thus applied to synergistically integrate AOD acquired from diversified sources. Consequently,  
256 creating the data tensor of AOD is of critical importance. In this study, the data tensor of AOD was  
257 constructed by incorporating not only observational AOD from both satellites and those inferred from  
258 *in situ* air quality indicators on the same date, but also historical AOD observations from MODIS  
259 instruments and part of data from the downscaled MERRA-2 AOD (denoted as AOD<sub>M2</sub> hereafter). The  
260 latter two were applied to provide knowledge of AOD distributions over space to guide the  
261 reconstruction of missing values in AOD<sub>Terra</sub>.

262 For the screening of historical observations resembling AOD<sub>Terra</sub> distribution on the given date,  
263 AOD<sub>M2</sub> was used in concert with AOD<sub>Terra</sub> and site-based AOD estimations to select similar imageries.  
264 Specifically, site-specific AOD estimations and 5% randomly selected AOD<sub>M2</sub> data were merged with  
265 valid AOD<sub>Terra</sub> to form a new image on each date, which was then used to find similar historical  
266 AOD<sub>Terra</sub> maps. To avoid the inclusion of imageries with distinct variation patterns, only those closely  
267 resembling AOD distribution in the composite image on the given date were selected in terms of their  
268 correlations and biases subject to a threshold of  $R > 0.7$  and  $RMSE < 0.2$ . Once sufficient historical  
269 imageries were obtained, the data tensor of AOD was constructed by compiling the observed AOD  
270 imageries on the given date with historical imageries to a three-dimension data array  $\mathbf{A} \in \mathbf{R}^{N_1 \times N_2 \times N_3}$   
271 (composed of  $N_3$  images with a size of  $N_1 \times N_2$ ). Data values of site-specific AOD estimations and 1%



272 randomly selected  $AOD_{M2}$  data were directly placed on grids where  $AOD_{Terra}$  values missed on each  
273 specific date. This greatly facilitated the reconstruction of missing AOD information over regions with  
274 tremendous data gaps in satellite observed AOD imageries given the presence of prior knowledge.  
275 More importantly, it significantly reduced the time required for convergence during the gap filling  
276 process.

### 277 3.2.2 Gap filling via tensor completion

278 Previous studies have well demonstrated the good performance of matrix decomposition  
279 methods such as empirical orthogonal function and singular value decomposition (SVD) for missing  
280 value imputation (Bai et al., 2020b; Beckers and Rixen, 2003; Folch-Fortuny et al., 2015). However,  
281 these methods can only work on two-dimension matrix mathematically, i.e., the matrix domain. To  
282 integrate spatial features of AOD revealed by datasets to generate a smooth AOD distribution with  
283 complete coverage, in this study, the HOSVD, a specific orthogonal Tucker decomposition, was  
284 applied. More detailed descriptions to HOSVD can be found in the literature such as Sun et al. (2021),  
285 Tucker (1966), Kolda and Bader (2009), Sidiropoulos et al. (2017), and Chen et al. (2014).

286 In Table 2, we provided a stepwise description of the algorithm used to fill data gaps in  $AOD_{Terra}$   
287 by integrating AOD features recognized in different imageries as the data tensor of AOD via HOSVD.  
288 To initiate the tensor decomposition, grids with missing values in the AOD tensor were first filled with  
289 the spatial average of valid AOD data in each individual image. Then, the AOD tensor was  
290 decomposed along each of three dimensions, while the dominant features in each dimension  
291 determined by the corresponding rank values were applied to reconstruct the data tensor. By gradually  
292 increasing the rank values and iteratively updating the initial filled values, the tensor can be  
293 reconstructed to better delineate AOD distribution over space after several iterations.

294 To confirm the convergence, a small portion of observational AOD values were randomly held out  
295 in advance, and the reconstructed values over these grids in each iteration were compared with these  
296 hold-out data till the difference between them lower than 0.01 (a threshold to determine convergence,  
297 a.k.a,  $\epsilon_1$  in Table 2). Meanwhile, to make the computational burden manageable, the study region  
298 (China in this study) was divided into 40 subregions (refer to Figure S2 for the spatial distribution of  
299 these subregions), and the tensor completion was then performed over each individual region. Finally,  
300 the reconstructed imageries were mosaiced to attain a national gap-free AOD map on each specific



301 date. During this step, a weighted average method was applied to solve the boundary effect when  
 302 mosaicking two adjacent maps, i.e., averaging the data value on each overlapped grid at the boundary  
 303 (50 km on the edge of subregion) as weighted by the distance to the edge. In the end, the mosaic  
 304 AOD<sub>Terra</sub> image was retained as the final gap-free AOD product.

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

<p><b>Input:</b> tensor <math>\mathbf{A} \in \mathbf{R}^{N_1 \times N_2 \times N_3}</math> with <math>\Omega = \{(i, j, k): A_{ijk} \text{ is observed}\}</math>, threshold <math>T_1, T_2</math></p> <p><b>Output:</b> reconstructed entries <math>\mathbf{A}' = \mathbf{A}^*(:, :, k^t) \in \mathbf{R}^{N_1 \times N_2}</math></p> <p>1: Initialize <math>A_{ijk}^* = \begin{cases} A_{ijk} &amp; (i, j, k) \in \Omega \\ \sum_i \sum_j A_{ijk} &amp; (i, j, k) \notin \Omega \end{cases}</math></p> <p>2: <b>for</b> <math>n_3 = N_3</math> to 1 <b>do</b></p> <p>3:     <math>n_1 = n_2 = 0</math></p> <p>4:     <b>while</b> <math>\varepsilon_1 &gt; T_1</math> <b>do</b></p> <p>5:         <math>n_1 = n_1 + 1, n_2 = n_2 + 1</math></p> <p>6:         Tucker Decomposition of <math>\mathbf{A}^*</math> with rank = <math>\{n_1, n_2, n_3\}</math>:                     <math>\mathbf{A}^* = \mathbf{S} \times_1 \mathbf{U}^{(n_1)} \times_2 \mathbf{U}^{(n_2)} \times_3 \mathbf{U}^{(n_3)}</math></p> <p>7:         <math>\varepsilon_1 = \arg \min_{\Omega} \frac{1}{2} \ \mathbf{A} - \mathbf{A}^*\ ^2</math></p> <p>8:         <math>\mathbf{A}_{\Omega}^* = \mathbf{A}_{\Omega}</math></p> <p>9:     <b>end while</b></p> <p>10:     <b>if</b> <math>\arg \min_{\Omega} \frac{1}{2} \ \mathbf{A} - \mathbf{A}^*\ ^2 &lt; T_2</math> <b>then</b></p> <p>11:         <b>break</b>;</p> <p>12:     <b>end if</b></p> <p>13: <b>end for</b></p>
--

### 306 3.3 PM concentration estimation

307 In this study, the random forest method was applied to establish regression models for PM<sub>2.5</sub> and  
 308 PM<sub>10</sub> concentration mapping. Ground measured PM<sub>2.5</sub> (or PM<sub>10</sub>) concentration data were used as the  
 309 learning target while AOD, aerosol components (AER<sub>comp</sub>), meteorological factors (MET), digital  
 310 elevation model (DEM), NDVI, land cover information (LC), and population were used as regressors.  
 311 The prediction model can be generally formulated as:

$$312 \quad \text{PM}_x = f(\text{AOD}, \text{AER}_{\text{comp}}, \text{MET}, \text{DEM}, \text{NDVI}, \text{POP}, \text{LC}, \text{month}) \quad (1)$$

313 where month is a categorical variable that was used to account for monthly varying relationships  
 314 between AOD and PM. For cross validation, PM<sub>2.5</sub> and PM<sub>10</sub> data from 10% of monitoring sites were  
 315 randomly held out to validate the predictive performance of each regression model. 500 regression  
 316 trees were used in each RF model, and each tree was grown on a bootstrap sample. The learning data



317 set was randomly divided into two parts during the training process, with 80% used as the training set  
318 while the rest 20% for testing. In order to guarantee a larger value of  $PM_{10}$  than  $PM_{2.5}$ ,  $PM_{2.5}$   
319 estimations from Eq. (1) were used as one predictor in addition to factors used to predict  $PM_{2.5}$  when  
320 estimating  $PM_{10}$  concentration. Such a model can also significantly improve the prediction accuracy  
321 of  $PM_{10}$  given the prior  $PM_{2.5}$  information.

### 322 **3.4 Point-surface data fusion**

323 Ground measured  $PM_{2.5}$  and  $PM_{10}$  concentration data were further fused with gridded PM  
324 estimations to enhance the data accuracy of PM data after 2014. Here, the well-known optimal  
325 interpolation (OI) method was applied to perform point-surface fusions between two different types  
326 datasets. Please refer to Bai et al. (2021) and Li et al. (2021) for a more detailed description of the OI  
327 method used to fuse PM concentration data. In this study, a modified scheme was developed to select  
328 neighboring observations. To avoid an isotropic interpolation effect, here we only used 30 ground  
329 observations with land cover, terrain and atmospheric conditions similar to those at the analyzed grid  
330 cell to estimate the innovation that should be assigned to the background value at the given grid. In  
331 other words, a similarity measure was first estimated between the analyzed grid cell and neighboring  
332 sites in terms of land cover, DEM, and atmospheric conditions. The 30 observations with similar  
333 background fields were then used in the OI procedure to correct possible bias in gridded PM  
334 estimations. Such a treatment can help exclude those observations with different ambient background,  
335 e.g., one site not far from the given grid but separated by a high mountain, thereby avoiding the possible  
336 propagation of antiphase corrections to data over adjacent grids.

## 337 **4 Results and discussion**

### 338 **4.1 Data accuracy of gap-free AOD in LGHAP**

339 Table 3 summarizes the data accuracy of gap-free AOD dataset generated in this study. For  
340 comparison, the data accuracy of each original AOD dataset was also assessed. Since *in situ* AOD  
341 measurements were not used as data input when reconstructing missing AOD information, thereby the  
342 gap-free AOD can be directly compared with *in situ* AOD measurements. As indicated, all these AOD  
343 datasets are in a good agreement with *in situ* AOD measurements. Generally, AODs from MODIS  
344 onboard Terra and Aqua have a similar data accuracy, which is also among the highest when



345 comparing with other datasets ( $R=0.95$  and  $RMSE=0.14$ ). AODs from AATSR show a comparable  
 346 accuracy with that of MODIS, but with a relatively low correlation with ground AOD measurements.  
 347 AODs from MISR, POLDER and VIIRS exhibit a similar bias level, with  $R$  varying from 0.80 to 0.92  
 348 and  $RMSE$  ranging from 0.22 to 0.29. In contrast,  $AOD_{M2}$  data have the poorest accuracy among these  
 349 eight gridded AOD datasets ( $R=0.77$  and  $RMSE=0.36$ ), even though AOD data from AERONET and  
 350 satellite observations like MODIS had been already assimilated. This indicates the presence of large  
 351 biases in  $AOD_{M2}$  and thus these  $AOD_{M2}$  data cannot be used solely to delineate AOD distributions  
 352 over space.

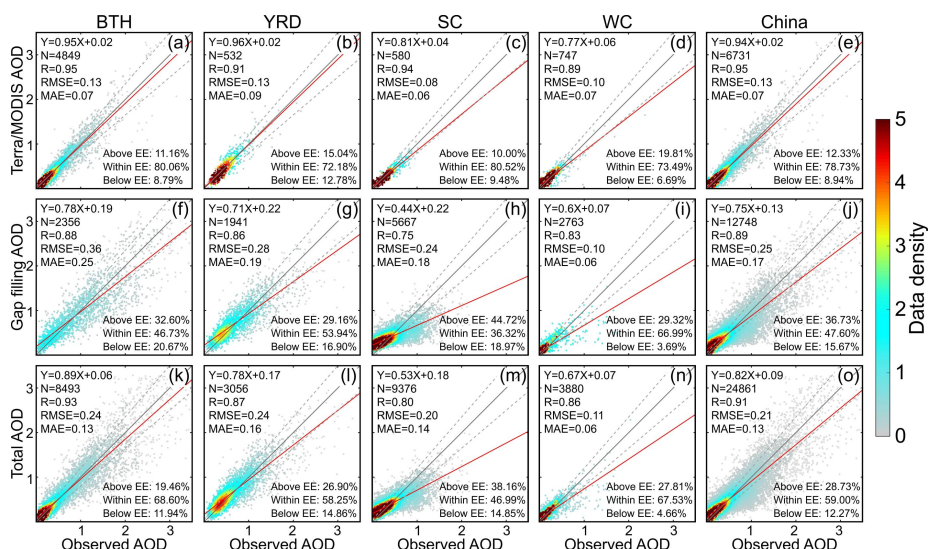
353 **Table 3.** Data accuracy of original and gap-free AOD datasets used and/or generated in this study. The  
 354 expected error (EE) was defined as  $\pm 0.05 + 0.15 \times AOD_{site}$ .

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

355  
 356 Compared to the first seven gridded AOD datasets, the LGHAP AOD dataset has an accuracy  
 357 slightly worse than the original MODIS AOD product but comparable to AODs from MISR, POLDER  
 358 and MERRA-2, with  $R$  of 0.91 and  $RMSE$  equaling to 0.21 compared to ground-based AOD  
 359 observations. Compared with AODs from MODIS, the gap-filled AOD appeared to overestimate  
 360 ground-based AOD observations, and this could be caused by the involvement of AODs from VIIRS  
 361 and POLDER as these two products of ground AOD observations significantly overestimated, which  
 362 can be indicated by the proportion of data pairs above the expected error (EE). On the other hand, such  
 363 significant underestimations in  $AOD_{M2}$  were not introduced to the LGHAP AOD as the former had a  
 364 below EE ratio of 32.97% which was only 12.27% in the latter. These results indicate that the gap-free  
 365 LGHAP AOD data are more likely to resemble AOD distributions revealed by satellite observations



366 rather than  $AOD_{M2}$ , justifying the advantages of involving multiple satellite AOD observations to help  
 367 reconstruct missing AOD values. Figure 2 further compares the data accuracy of  $AOD_{Terra}$  and the  
 368 reconstructed data over different regions of China. It is indicative that the purely reconstructed data  
 369 have an accuracy ( $R=0.88$  and  $RMSE=0.26$ ) lower than the original  $AOD_{Terra}$  ( $R=0.95$  and  
 370  $RMSE=0.13$ ) across China, especially in South China where the reconstructed data were significantly  
 371 underestimated the ground-based AOD observations. Possible reasons for the relatively poor accuracy  
 372 of AOD reconstructions in this region could be attributed to extensive data gaps over there due to  
 373 frequent clouds (refer to Figure S3 for the distribution of mean data integrity of  $AOD_{Terra}$  during 2000–  
 374 2020), which significantly limit the learning capacity in space and temporal domain during the tensor  
 375 completion process. In other words, limited observations in satellite imageries greatly reduced the  
 376 learning performance from the sparse tensor. Even though, the purely reconstructed data exhibit a bias  
 377 level comparable to AOD retrievals from several satellite instruments, e.g., MISR, VIIRS, and  
 378 POLDER. This demonstrates the good performance of the proposed tensor completion method in  
 379 reconstructing missing AOD information. By combining the reconstructed data with original AOD  
 380 observations from Terra/MODIS, we obtained a 21-year-long spatially complete high-resolution  
 381 (daily/1-km) AOD product with satisfying accuracy ( $R=0.91$  and  $RMSE=0.21$ ).



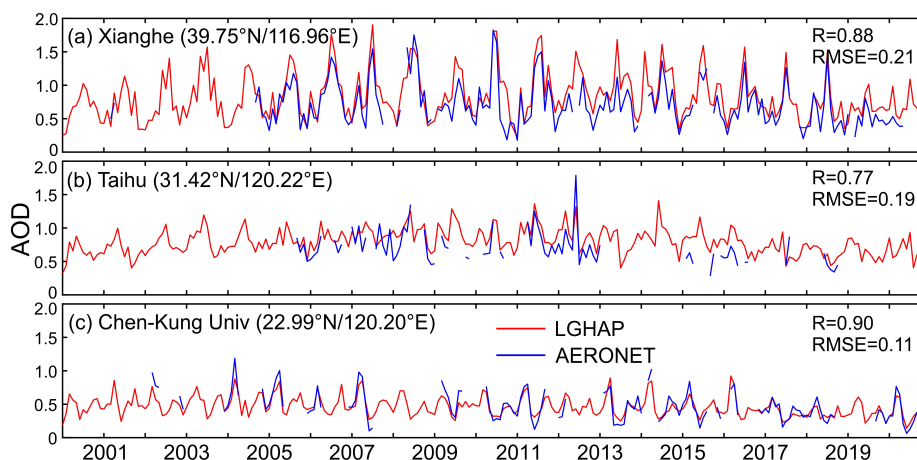
382  
 383 **Figure 2.** Scatter plots between ground observed and satellite-based AOD data in different regions of  
 384 China. (a–e) original Terra/MODIS AOD, (f–j) reconstructed AOD only, and (k–o) both original and





385 reconstructed data combined. BTH, YRD, SC, and WC refers to regions of Beijing-Tianjin-Hebei,  
386 Yangtze River Delta, South China, and West China, respectively.

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

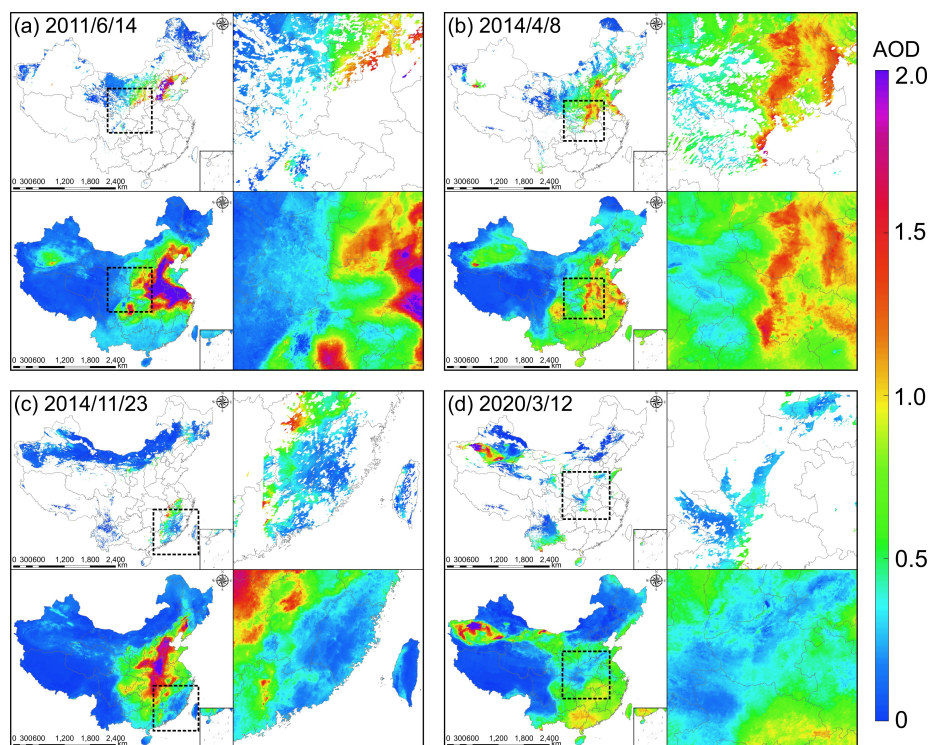


401  
402 **Figure 3.** Comparison of monthly AOD time series from LGHAP and AERONET at three different  
403 stations in China. Latitude and longitude information of each site was given in brackets.

404 Since the final gap-free AOD product was generated mainly by integrating a set of data tensor  
405 of gridded AOD with *in situ* AOD estimations, the relative contribution of each product to the final



406 gap-free dataset is worth being investigated. In this study, a data coverage ratio weighted nonlinear  
407 correlation coefficient was proposed to examine the relative contribution of each gridded product to  
408 the LGHAP AOD dataset. The nonlinear correlation coefficient was used to assess the mutual  
409 information between two variables (Sun et al., 2021; Wang et al., 2005), while the data coverage ratio  
410 was multiplied to indicate the overall contribution of one product to the final fused dataset (refer to  
411 Text S2 for the definition of this indicator). As shown in Figure 5, the relative contribution of each  
412 gridded product varied with time and the input data sources. In the early two years (2000–2001), the  
413 AOD distribution in gap-free imageries was determined largely by AOD<sub>Terra</sub> (81%), whereas this ratio  
414 decreased to about 30% when many other products were involved, especially AOD from Aqua and  
415 PARASOL. With the advent of VIIRS and the loss of PARASOL after 2012, the relative contribution  
416 changed drastically as AOD from MODIS and VIIRS played the dominant roles in reconstructing  
417 AOD distribution. Note the relative contribution of AOD<sub>M2</sub> remained lower than 10%, indicative of  
418 the greater importance of satellite observations in generating the LGHAP AOD product.

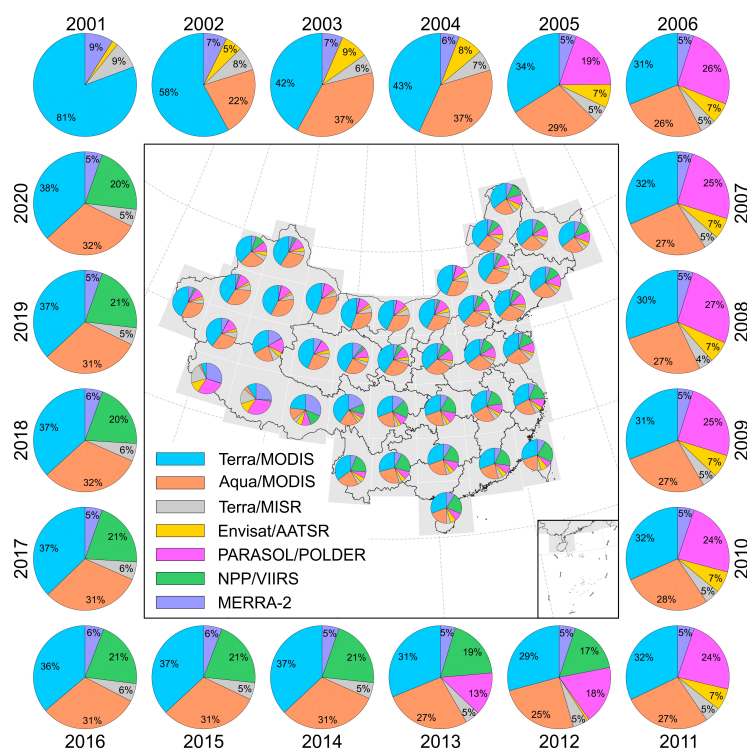


419  
420 **Figure 4.** Spatial patterns of the reconstructed AOD under different baseline AOD coverage ratios. In  
421 each sub-diagram, the upper panel presents the original AOD distribution from Terra/MODIS while



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

424 With respect to the temporally averaged contribution in each subregion, it shows that the  
425 relative contribution of each product also varied significantly across regions. Generally, AOD from  
426 MODIS aboard Terra and Aqua played the most important role (>60%) in generating the LGHAP  
427 AOD product, except over the southwest part of the country (Tibet plateau) where AOD<sub>M2</sub> contributed  
428 most. This is reasonable since data gaps are abnormally high in satellite observations over this region  
429 because of the vast and long-lasting snow cover (refer to Figure S3 for the data integrity distribution).  
430 Consequently, AOD<sub>M2</sub> would play an important role in reconstructing AOD distribution over such  
431 regions. Overall, the results shown here clearly highlight the success of big data analytics in generating  
432 the LGHAP AOD dataset via integrative efforts from diversified data sources.



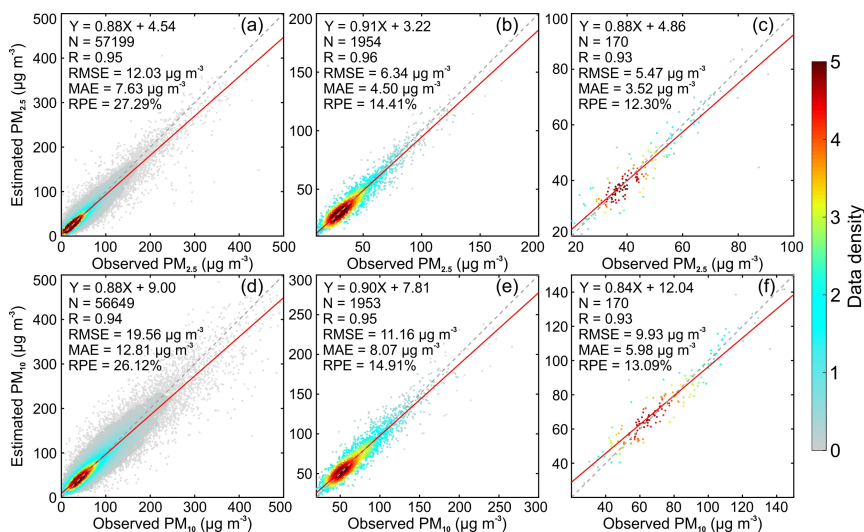
433  
434 **Figure 5.** Spatiotemporal variations of the relative contribution of each gridded AOD product to the  
435 generation of LGHAP AOD dataset. The relative contribution was estimated as the data coverage ratio  
436 weighted nonlinear correlation coefficient (please refer to Text S2 in the supplementary information  
437 for the arithmetic theory to calculate this measure). The annual mean shown outside is the national



438 averaged contribution in each individual year while the regional mean shown on the map was averaged  
 439 over the past 21-year in each subregion.

#### 440 4.2 Data accuracy of PM<sub>2.5</sub> and PM<sub>10</sub> estimations

441 By taking advantage of the gap-filled AOD, daily 1-km resolution PM<sub>2.5</sub> and PM<sub>10</sub> concentration  
 442 data in China were estimated via an ensemble learning approach. Figure S4 shows the sample-based  
 443 cross validation accuracy of two prediction models. It shows that the original daily PM<sub>2.5</sub> prediction  
 444 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 accuracy is  
 445 comparable to our previous study (Bai et al., 2019a), but slightly worse than those reported in some  
 446 recent studies (Table 4). In contrast, PM<sub>10</sub> had a much higher prediction accuracy, with R<sup>2</sup> of 0.90 and  
 447 RMSE of 21.06 μg m<sup>-3</sup> for the daily product. This good performance should be attributed to the  
 448 involvement of PM<sub>2.5</sub> estimations as a predictor in the PM<sub>10</sub> prediction model. Figure 6 shows the site-  
 449 specific (held-out in advance) validation accuracy of daily, monthly, and annual mean PM<sub>2.5</sub> and PM<sub>10</sub>  
 450 concentration in LGHAP. As shown, the site-specific validation results indicated that the final full-  
 451 coverage (gap free) daily PM<sub>2.5</sub> and PM<sub>10</sub> concentration data are in a good agreement with ground-  
 452 based measurements, with R of 0.95 and RMSE of 12.03 μg m<sup>-3</sup> for PM<sub>2.5</sub> while R of 0.94 and RMSE  
 453 of 19.56 μg m<sup>-3</sup> for PM<sub>10</sub>. Overall, PM data in LGHAP are not only spatially complete with a finer  
 454 resolution but have a comparable accuracy with previous studies.



455



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

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

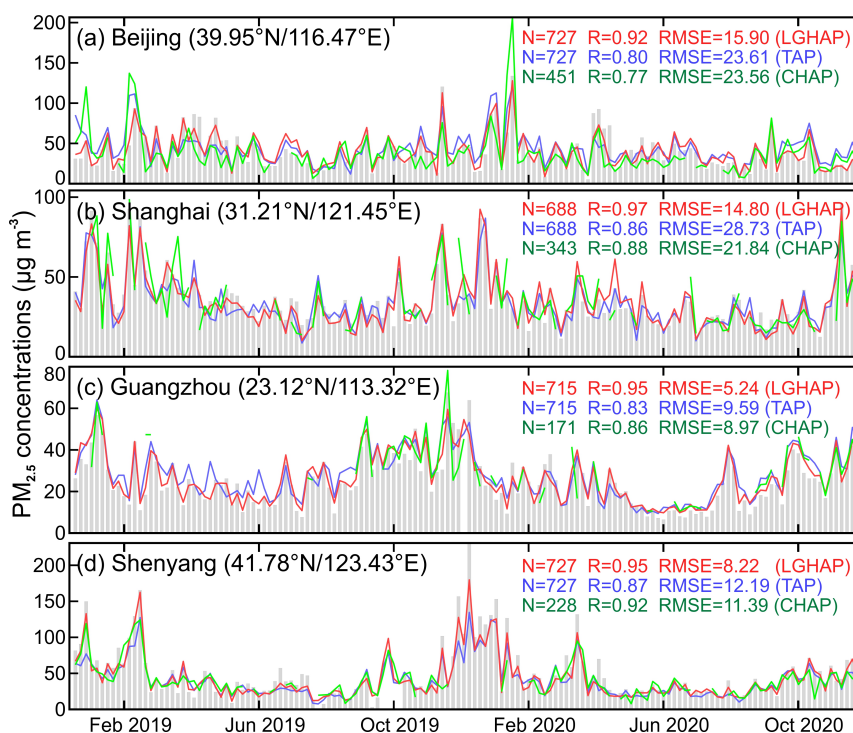
Source	Gap-free	Resolution	Time range	$R^2$	RMSE ( $\mu\text{g m}^{-3}$ )
Wei et al. (2021)	No	1 km	2000~2018	0.86~0.90	10.09~18.39
Geng et al. (2021)	Yes	10 km	2000~2021	0.80~0.88	13.90~22.10
Xue et al. (2019)	Yes	10 km	2000~2016	0.61	27.80
Chen et al. (2018)	No	10 km	2005~2016	0.83	28.10
Lyu et al. (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 $PM_{2.5}$	Yes	1 km	2000~2020	0.90	12.03

461

462 Figure 7 presents a two-year-long comparison of  $PM_{2.5}$  concentration time series from LGHAP  
463 and two other open access datasets with  $PM_{2.5}$  measurements sampled at four United States Embassy  
464 in China. Since this ground-based dataset has been seldomly noticed and used, it can be applied as an  
465 independent dataset to fairly evaluate the accuracy of these three machine-learned  $PM_{2.5}$  estimations.  
466 As shown, all these three datasets well reconstructed temporal variations of  $PM_{2.5}$  from 2019 to 2020.  
467 Temporally, LGHAP and TAP are continuous while CHAP suffers from significant data gaps because  
468 no gap filling method was applied when generating the dataset. Compared with the other two datasets,  
469 LGHAP  $PM_{2.5}$  data had a better agreement with ground-based  $PM_{2.5}$  measurements. This high  
470 accuracy could be partially due to the fusion of *in situ*  $PM_{2.5}$  data measured at adjacent sites via the OI  
471 method. Figure S5 compares  $PM_{2.5}$  time series from LGHAP with  $PM_{2.5}$  measurements sampled at  
472 five United States Embassy in China. It is indicative that historical  $PM_{2.5}$  variations over these five  
473 cities were well reconstructed in LGHAP, even over years before 2014 at which  $PM_{2.5}$  measurements  
474 from state-control monitoring sites were not available. Note  $PM_{2.5}$  estimations appeared to  
475 significantly underestimate  $PM_{2.5}$  concentration sampled at the Embassy in Beijing before 2013.



476 Considering the reconstructed AOD time series agreed well with AERONET AOD in Beijing (Figure  
477 3a), and the model performed well in predicting historical PM<sub>2.5</sub> in Shanghai during the synchronous  
478 time period (Figure S5b), we are more willing to attribute this issue to significant PM<sub>2.5</sub>  
479 overestimations by the US Embassy during that period. Overall, these independent validation results  
480 collectively indicate a good accuracy of PM<sub>2.5</sub> in LGHAP dataset.



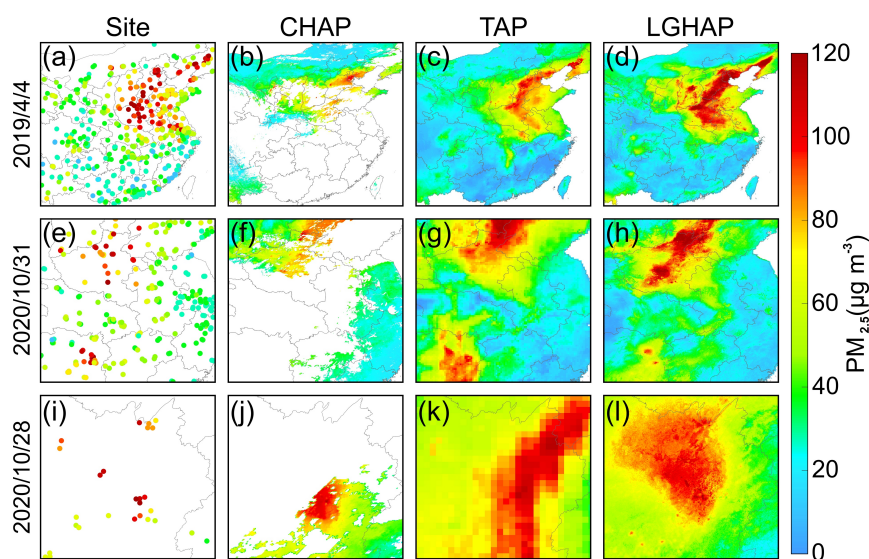
481  
482 **Figure 7.** Comparison of PM<sub>2.5</sub> concentration time series between LGHAP (red line) and two open  
483 datasets (blue: TAP, green: CHAP). Here, hourly PM<sub>2.5</sub> concentrations measured by four United States  
484 Embassy in China from 2019 to 2020 (grey bar) were used as an independent PM<sub>2.5</sub> dataset to validate  
485 these three daily products. CHAP and TAP are two open access datasets providing PM<sub>2.5</sub>  
486 concentration that were created by Wei et al. (2021) and Geng et al. (2021) respectively.

487

488 In Figure 8 we compared the spatial distribution of PM<sub>2.5</sub> that was reconstructed by different  
489 datasets. Compared to LGHAP and TAP, PM<sub>2.5</sub> data from CHAP are not gap free since the spatial  
490 coverage is determined by the AOD data coverage in the MAIAC product. Compared to TAP, LGHAP  
491 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



492 space with more details. Overall, LGHAP has a better performance in reconstructing  $PM_{2.5}$  spatial  
 493 distributions than the other two datasets. Reasons could be attributed to the following two aspects.  
 494 Firstly, in situ  $PM_{2.5}$  measurements were fused with gridded  $PM_{2.5}$  estimations using the OI method  
 495 when generating the final  $PM_{2.5}$  product in LGHAP. This can help correct modeling biases in original  
 496  $PM_{2.5}$  estimations. Secondly, a set of satellite-based AOD retrievals were incorporated when  
 497 generating the full-coverage AOD product, which greatly helps reduce large biases in numerical AOD  
 498 simulations, yielding more accurate  $PM_{2.5}$  estimations in turn. This also highlights the great advantages  
 499 of using big data analytics methods to advance air pollution assessment.

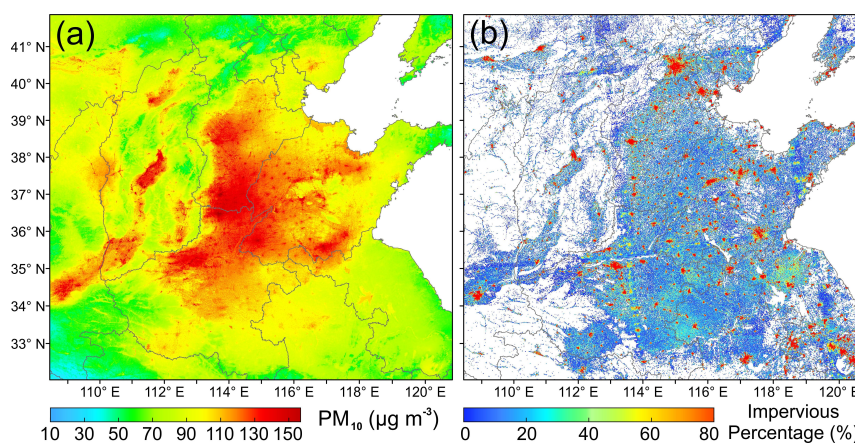


500  
 501 **Figure 8.** Comparison of  $PM_{2.5}$  distribution reconstructed by different  $PM_{2.5}$  concentration datasets.  
 502 From the left to right, it shows in situ  $PM_{2.5}$  concentration measurements, CHAP, TAP, and LGHAP,  
 503 respectively.

504  
 505 To illustrate the fine resolution of LGHAP dataset, we compared the annual mean  $PM_{10}$   
 506 concentration in 2019 with the proportion of impervious surface that was derived from 30-m resolution  
 507 land cover dataset in eastern China. As shown in Figure 9, the finer resolution of LGHAP dataset  
 508 enables us to easily recognize the “hot spot” regions with high  $PM_{10}$  loading. By referring to the  
 509 impervious surface distribution on the right, we found that these hot spots are mainly over cities and  
 510 towns, indicative of the presence of pollution island in urban regions. Owing to the involvement of  
 511 such high-resolution datasets, the spatial details of  $PM_{2.5}$  and  $PM_{10}$  can be then well recognized in



512 LGHAP. The finer spatial resolution advantage of the LGHAP dataset can be also demonstrated by  
513 comparisons of spatial distribution of annual mean  $PM_{2.5}$  concentration that was revealed by four  
514 different datasets shown in Figure S6.



515  
516 **Figure 9.** Comparison of annual mean  $PM_{10}$  concentration with the proportion of areas covered by  
517 impervious surface in eastern China.

### 518 4.3 Long-term trends of haze pollution in China from 2000 to 2020

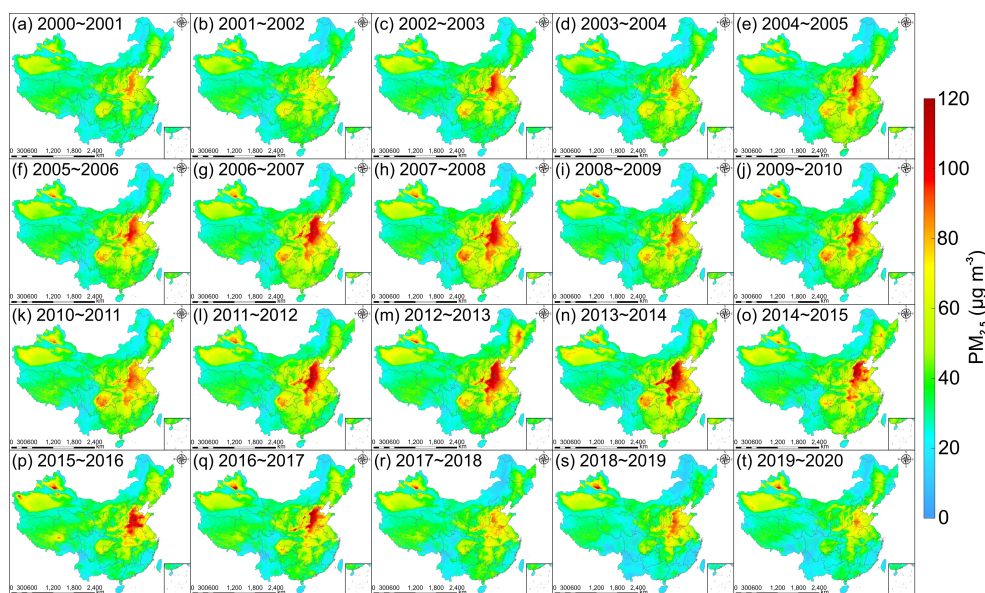
519 The aerosol pollution trends in China can be better examined by taking advantage of LGHAP  
520 dataset given long temporal coverage, gap free and high-resolution superiorities. Severe haze  
521 pollutions such as  $PM_{2.5}$  are oftentimes observed during the wintertime (September to February). In  
522 this study, we first calculated wintertime mean  $PM_{2.5}$  concentration in China from 2000 to 2020. As  
523 shown in Figure 10, severe wintertime haze pollution events were mainly observed in North China,  
524 especially over the adjacent region in Hebei-Shandong-Henan provinces. In addition, Sichuan basin  
525 and Fenwei plain also suffered from severe haze pollution during the wintertime. Temporally, severe  
526 haze pollution events occurred mainly from the late 2002 to early 2017, which were significantly  
527 reduced after 2017. Similar pattern can be also inferred from  $PM_{10}$  concentration distributions shown  
528 in Figure S7.

529 Figure 11 shows the temporal variations of the proportion of land areas covered by  $PM_{2.5}$   
530 concentration exceeding  $35 \mu g m^{-3}$  (the national ambient air quality standard for 24-hour  $PM_{2.5}$   
531 concentration given in GB 3095-2012). As shown in Figure 11a, severe  $PM_{2.5}$  pollution occurred  
532 mainly during the wintertime in China, as more than one-third land areas (indicated by the blue lines)





533 were exposed to hazardous  $PM_{2.5}$  pollutants. Meanwhile, an apparent inflection was observed in 2007,  
 534 after which the number of episode days decreased drastically at more than one-third land area covered  
 535 by  $PM_{2.5}$  concentration exceeding  $35 \mu g m^{-3}$ . According to the proportion of land area covered with  
 536 annual mean  $PM_{2.5}$  concentration greater than  $35 \mu g m^{-3}$ , the variation of haze pollution in China can  
 537 be generally divided into three different periods during the past two-decades (Figure 11b). As indicated,  
 538 an increasing trend was observed from 2000 to 2007, during which land areas covered by  $PM_{2.5}$   
 539 concentration greater than  $35 \mu g m^{-3}$  had increased to near 40% at a pace of  $1.04\% a^{-1}$ . The second  
 540 period was from 2008 to 2013, during which the land area coverage ratio decreased at a rate of  $-0.21\%$   
 541  $a^{-1}$ . The third period started from 2014, after which the land area covered with  $PM_{2.5}$  concentration  
 542 more than  $35 \mu g m^{-3}$  had decreased drastically, at a rate of  $-2.23\% a^{-1}$ .

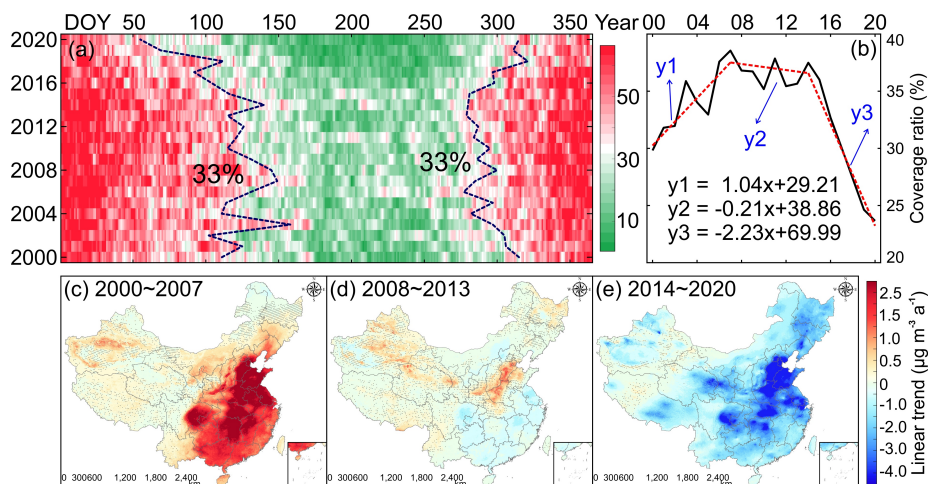


543  
 544 **Figure 10.** Spatial distribution of wintertime (September to February) averaged  $PM_{2.5}$  concentration  
 545 from LGHAP during 2000 to 2020 in China.  
 546

547 Figure 11c–e presents the linear trend of  $PM_{2.5}$  concentration during these three specific periods,  
 548 from which we observed that significant  $PM_{2.5}$  variations occurred mainly over eastern parts of the  
 549 country where resides two-thirds of the population. A near ubiquitous  $PM_{2.5}$  increasing trend was  
 550 observed during 2000–2007, with significant increase ( $>1.0 \mu g m^{-3} a^{-1}$ ) mainly observed in eastern  
 551 China. During the second period,  $PM_{2.5}$  concentration over most regions shows a small decreasing



552 trend except in the Ji-Lu-Yu region where an increasing trend was still observed. Apparent decreasing  
553 trend was observed over most parts of the country after 2014, indicative of significant reductions in  
554 PM<sub>2.5</sub> loading across China. This trend distribution is in line with our previous study that was derived  
555 using the annual mean PM<sub>2.5</sub> concentration dataset generated by the Dalhousie University (Bai et al.,  
556 2019b). However, differences were still observed in terms of the regions where significant decreasing  
557 trends were present. Most significant decreasing trends were mainly observed in Sichuan basin and  
558 Pearl River Delta in the previous study. However, regions with drastic PM<sub>2.5</sub> decrease were found  
559 mainly in the North China where severe haze pollution events were oftentimes reported. Similar  
560 variation patterns can be also inferred from PM<sub>10</sub> (Figure S8) and AOD (Figure S9). Overall, the  
561 LGHAP dataset provides us a gridded perspective to better examine long-term variations of haze  
562 pollution in China during the past two decades.



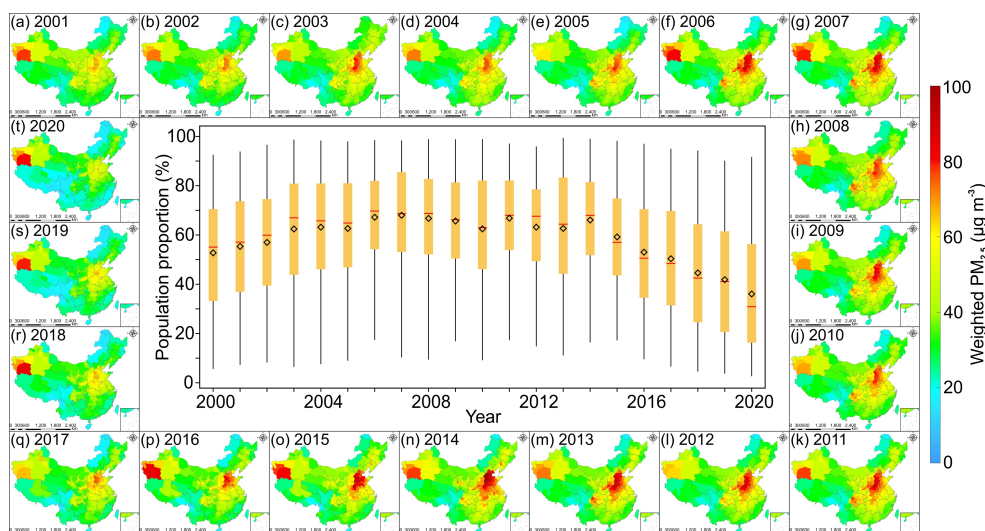
563  
564 **Figure 11.** Temporal variations of the proportion of land areas covered with PM<sub>2.5</sub> concentration  
565 exceeding 35 µg m<sup>-3</sup> and PM<sub>2.5</sub> trends during three different periods. (a) Temporal variations of the  
566 land coverage ratio with daily PM<sub>2.5</sub> concentration exceeding 35 µg m<sup>-3</sup> from 2000 to 2000. (b) same  
567 as (a) but for annual mean PM<sub>2.5</sub> concentration. (c–e) PM<sub>2.5</sub> trends during periods of 2000–2007, 2008–  
568 2013, and 2014–2020. The dotted regions imply trend estimations are statistically insignificant at the  
569 95% confidence interval.

#### 570 4.4 Population exposure to PM<sub>2.5</sub> pollution in China



571 By taking advantage of fine resolution LGHAP PM<sub>2.5</sub> concentration and gridded population data,  
572 population exposure to PM<sub>2.5</sub> pollution across China over the past two decades were estimated. Figure  
573 12 shows the spatial distribution of population weighted PM<sub>2.5</sub> concentration and the proportion of  
574 population exposed to PM<sub>2.5</sub> concentration greater than 35  $\mu\text{g m}^{-3}$ . As shown, spatial distribution of  
575 population weighted PM<sub>2.5</sub> concentration resembles the spatial pattern of annual mean PM<sub>2.5</sub>  
576 concentration, with high values observed mainly in eastern and central China as well as northwest  
577 China. Nonetheless, PM<sub>2.5</sub> sources in these two areas could be different. In northwest China, natural  
578 emissions could be the dominant source since very limited population resides there. In contrast, most  
579 population lives in eastern and central China with highly developed economy, and anthropogenic  
580 emissions thus might play more important roles in PM<sub>2.5</sub> formation (Xin et al., 2015; Yang et al., 2011).  
581 In regard to the proportion of population exposed to the ambient with PM<sub>2.5</sub> concentration greater than  
582 35  $\mu\text{g m}^{-3}$ , we observed that the annual mean population ratio exposure to unhealthy PM<sub>2.5</sub> increased  
583 gradually from 50.60% in 2000 to 65.72% in 2007. During 2007–2014, the ratio varied with small  
584 changes (<5%), whereas a drastic decline was observed after 2014, with the annual mean proportion  
585 of population exposed to unhealthy PM<sub>2.5</sub> was reduced from 63.81% in 2014 to 34.03% in 2020, even  
586 though the total population was increased from 1.37 billion to 1.41 billion during the synchronous  
587 period. Nonetheless, more than one-third population was still exposed to unhealthy PM<sub>2.5</sub>, highlighting  
588 the requirement of further emission reduction actions to manage haze pollutions in China.

589





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

## 595 **5 Data availability**

596 The LGHAP dataset, consisting of gap free AOD,  $PM_{2.5}$ , and  $PM_{10}$  concentration with daily 1-  
597 km resolution from 2000 to 2020, are all publicly accessible. The daily map was provided in the  
598 NetCDF format and data in each individual year were archived in a zip file. For AOD, the dataset has  
599 a disk storage size of near 27 GB in total, which can be found at <https://doi.org/10.5281/zenodo.5652257>  
600 (Bai et al., 2021a).  $PM_{2.5}$  (38 GB) and  $PM_{10}$  (48 GB) concentration data can be acquired from  
601 <https://doi.org/10.5281/zenodo.5652265> (Bai et al., 2021b) and <https://doi.org/10.5281/zenodo.5652263> (Bai  
602 et al., 2021c), respectively. Additionally, monthly and annual mean datasets were also provided, which  
603 can be acquired from <https://doi.org/10.5281/zenodo.5655797> (Bai et al., 2021d) and  
604 <https://doi.org/10.5281/zenodo.5655807> (Bai et al., 2021e), respectively.

## 605 **6 Conclusion**

606 In this study, a big data analytics method was developed for generating a LGHAP dataset to  
607 advance research in earth system science and environment management. With integrative efforts of  
608 fusing AOD features extracted from a set of AOD data tensors and knowledge transfer in statistical  
609 data mining from diverse air quality indicators, a LGHAP aerosol dataset providing 21-year-long  
610 (2000–2020) gap-free AOD,  $PM_{2.5}$ , and  $PM_{10}$  concentration data with daily 1-km resolution in China,  
611 was generated. Gap-filled AOD imageries were firstly generated by reconstructing AOD distribution  
612 in  $AOD_{Terra}$  via fusing AOD features recognized from diversified satellites and numerical models as  
613 well as in situ data through tensor completion. Compared to ground-based AOD measurements, the  
614 gap-filled AOD data exhibit a satisfying prediction accuracy and good performance in delineating  
615 AOD variations over space and time. To our knowledge, this is the first thrust of generating long-term  
616 high-resolution AOD dataset with gap free nature in China.



617 PM<sub>2.5</sub> and PM<sub>10</sub> concentration data were then estimated using an ensemble learning approach by  
618 taking advantage of the generated gap-free AOD imageries. Ground validation results also indicate  
619 good accuracies of these two gridded products, showing a comparable bias level with many previous  
620 studies. Compared with other open access daily PM<sub>2.5</sub> concentration datasets, the LGHAP PM<sub>2.5</sub>  
621 dataset performs well due to the vantage of having gap free and fine resolution products. With this gap  
622 free and high-resolution dataset, the long-term variation trend of haze pollution in China over the past  
623 two decades was examined, and apparent inflections were observed in 2007 and 2014, at which PM<sub>2.5</sub>  
624 concentration was found to turn from an increasing path to decreasing in 2007 with a more drastic  
625 decline observed starting from 2014. Moreover, the LGHAP dataset provides us a gridded perspective  
626 to assess two-decade long population exposure to PM<sub>2.5</sub> pollution in China. In spite of a drastic decline  
627 in population exposure, there are still more than one-third population exposed to unhealthy PM<sub>2.5</sub>  
628 pollutants, highlighting the requirement of long-lasting actions to continue PM<sub>2.5</sub> related emission  
629 reduction.

630 Overall, these three gridded LGHAP aerosol products provide a long-term perspective on aerosol  
631 changes over different regions of China, and users are encouraged to use the LGHAP dataset to assess  
632 aerosol impacts on public health, air quality, climate, and ecosystem. The dataset has been publicly  
633 released online and is freely accessible via the links provided above. In addition to the LGHAP dataset,  
634 Python, Matlab, R, and IDL codes that can be used to read and visualize these data were provided as  
635 well. Global scale dataset is on the track and will be released to the public soon.

#### 636 **Author contributions**

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

#### 642 **Competing interests**

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



644 **Acknowledgments**

645 The authors are grateful to all organizations and groups for providing essential datasets that were used  
646 in this study. The MAIAC AOD was acquired from <https://lpdaac.usgs.gov/products/mcd19a2v006/>.  
647 The MISR AOD was acquired from <https://asdc.larc.nasa.gov/project/MISR>. The VIIRS AOD was  
648 acquired from [https://earthdata.nasa.gov/earth-observation-data/near-real-time/download-nrt-](https://earthdata.nasa.gov/earth-observation-data/near-real-time/download-nrt-data/viirs-nrt)  
649 [data/viirs-nrt](https://earthdata.nasa.gov/earth-observation-data/near-real-time/download-nrt-data/viirs-nrt). The AATSR AOD was acquired from <https://climate.esa.int/en/projects/aerosol/data/>.  
650 The POLDER AOD was acquired from <https://www.grasp-open.com/products/polder-data-release/>.  
651 The aerosol diagnostics including AOD and aerosol components from MERRA-2 were acquired from  
652 [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  
653 AERONET was acquired from [https://aeronet.gsfc.nasa.gov/new\\_web/aerosols.html](https://aeronet.gsfc.nasa.gov/new_web/aerosols.html). Meteorological  
654 factors were retrieved from the latest ERA-5 reanalysis and can be reached at  
655 <https://cds.climate.copernicus.eu/cdsapp#!/dataset/reanalysis-era5-single-levels?tab=overview>.  
656 Atmospheric visibility data were acquired from the national meteorological information center at  
657 <http://data.cma.cn/en>. Ground-based air pollutants concentration was acquired from  
658 <https://air.cnemc.cn:18007/>. Gridded Population data were acquired from <https://www.worldpop.org/>  
659 while DEM was acquired from <https://www.resdc.cn/>. Monthly NDVI data were acquired from  
660 <https://lpdaac.usgs.gov/products/mod13a3v061/>. Land cover data were acquired from  
661 <http://www.globallandcover.com/defaults.html?src=/Scripts/map/defaults/browse.html&head=brows>  
662 [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>.

663 **Financial support**

664 This study was supported by the National Natural Science Foundation of China (grants 42171309 and  
665 41701413), and the Shanghai Committee of Science and Technology (grant 20ZR1415900).

666



667 **References**

668 Bai, K., Chang, N.-B. and Chen, C.-F.: Spectral Information Adaptation and Synthesis Scheme  
669 for Merging Cross-Mission Ocean Color Reflectance Observations From MODIS and VIIRS, IEEE  
670 Trans. Geosci. Remote Sens., 54(1), 311–329, doi:10.1109/TGRS.2015.2456906, 2016.

671 Bai, K., Li, K., Chang, N.-B. and Gao, W.: Advancing the prediction accuracy of satellite-based  
672 PM<sub>2.5</sub> concentration mapping: A perspective of data mining through in situ PM<sub>2.5</sub> measurements,  
673 Environ. Pollut., 254, 113047, doi:10.1016/j.envpol.2019.113047, 2019a.

674 Bai, K., Ma, M., Chang, N.-B. and Gao, W.: Spatiotemporal trend analysis for fine particulate  
675 matter concentrations in China using high-resolution satellite-derived and ground-measured PM<sub>2.5</sub>  
676 data, J. Environ. Manage., 233, 530–542, doi:10.1016/j.jenvman.2018.12.071, 2019b.

677 Bai, K., Li, K., Wu, C., Chang, N.-B. and Guo, J.: A homogenized daily in situ PM<sub>2.5</sub>  
678 concentration dataset from the national air quality monitoring network in China, Earth Syst. Sci. Data,  
679 12(4), 3067–3080, doi:10.5194/essd-12-3067-2020, 2020a.

680 Bai, K., Li, K., Guo, J., Yang, Y. and Chang, N.-B.: Filling the gaps of in situ hourly PM<sub>2.5</sub>  
681 concentration data with the aid of empirical orthogonal function analysis constrained by diurnal cycles,  
682 Atmos. Meas. Tech., 13(3), 1213–1226, doi:10.5194/amt-13-1213-2020, 2020b.

683 Bai, K., Li, K., Guo, J. and Chang, N.-B.: Multiscale and multisource data fusion for full-coverage  
684 PM<sub>2.5</sub> concentration mapping: Can spatial pattern recognition come with modeling accuracy?, ISPRS  
685 J. Photogramm. Remote Sens., 2021. in revision

686 Bai, K., Li, K. Tan, Z., Han, D., and Guo, J.: Daily 1-km gap-free AOD grids in China, v1 (2000–  
687 2020) [data set], <https://doi.org/10.5281/zenodo.5652257>, 2021a.

688 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–  
689 2020) [data set], <https://doi.org/10.5281/zenodo.5652265>, 2021b.

690 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–  
691 2020) [data set], <https://doi.org/10.5281/zenodo.5652263>, 2021c.

692 Bai, K., Li, K. Tan, Z., Han, D., and Guo, J.: Monthly averaged 1-km gap-free AOD, PM<sub>2.5</sub> and  
693 PM<sub>10</sub> grids in China, v1 (2000–2020) [data set], <https://doi.org/10.5281/zenodo.5655797>, 2021d.

694 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>  
695 grids in China, v1 (2000–2020) [data set], <https://doi.org/10.5281/zenodo.5655807>, 2021e.

696 Beckers, J. M. and Rixen, M.: EOF Calculations and Data Filling from Incomplete Oceanographic



- 697 Datasets, *J. Atmos. Ocean. Technol.*, 20(12), 1839–1856, doi:10.1175/1520-  
698 0426(2003)020<1839:ECADFF>2.0.CO;2, 2003.
- 699 Bi, J., Belle, J. H., Wang, Y., Lyapustin, A. I., Wildani, A. and Liu, Y.: Impacts of snow and cloud  
700 covers on satellite-derived PM<sub>2.5</sub> levels, *Remote Sens. Environ.*, 221(October), 665–674,  
701 doi:10.1016/j.rse.2018.12.002, 2018.
- 702 Chang, N.-B., Bai, K. and Chen, C.-F.: Smart Information Reconstruction via Time-Space-  
703 Spectrum Continuum for Cloud Removal in Satellite Images, *IEEE J. Sel. Top. Appl. Earth Obs.*  
704 *Remote Sens.*, 8(5), 1898–1912, doi:10.1109/JSTARS.2015.2400636, 2015.
- 705 Che, H., Yang, L., Liu, C., Xia, X., Wang, Y., Wang, H., Wang, H., Lu, X. and Zhang, X.: Long-  
706 term validation of MODIS C6 and C6.1 Dark Target aerosol products over China using CARSNET  
707 and AERONET, *Chemosphere*, 236, 124268, doi:10.1016/j.chemosphere.2019.06.238, 2019.
- 708 Chen, G., Li, S., Knibbs, L. D., Hamm, N. A. S., Cao, W., Li, T., Guo, J., Ren, H., Abramson, M.  
709 J. and Guo, Y.: A machine learning method to estimate PM<sub>2.5</sub> concentrations across China with  
710 remote sensing, meteorological and land use information, *Sci. Total Environ.*, 636, 52–60,  
711 doi:10.1016/j.scitotenv.2018.04.251, 2018.
- 712 Di, Q., Amini, H., Shi, L., Kloog, I., Silvern, R., Kelly, J., Sabath, M. B., Choirat, C., Koutrakis,  
713 P., Lyapustin, A., Wang, Y., Mickley, L. J. and Schwartz, J.: An ensemble-based model of PM<sub>2.5</sub>  
714 concentration across the contiguous United States with high spatiotemporal resolution, *Environ. Int.*,  
715 130, 104909, doi:10.1016/j.envint.2019.104909, 2019.
- 716 van Donkelaar, A., Martin, R. V., Brauer, M., Kahn, R., Levy, R., Verduzco, C. and Villeneuve,  
717 P. J.: Global Estimates of Ambient Fine Particulate Matter Concentrations from Satellite-Based  
718 Aerosol Optical Depth: Development and Application, *Environ. Health Perspect.*, 118(6), 847–855,  
719 doi:10.1289/ehp.0901623, 2010.
- 720 van Donkelaar, A., Martin, R. V., Brauer, M., Hsu, N. C., Kahn, R. A., Levy, R. C., Lyapustin,  
721 A., Sayer, A. M. and Winker, D. M.: Global Estimates of Fine Particulate Matter using a Combined  
722 Geophysical-Statistical Method with Information from Satellites, Models, and Monitors, *Environ. Sci.*  
723 *Technol.*, 50(7), 3762–3772, doi:10.1021/acs.est.5b05833, 2016.
- 724 Fang, X., Zou, B., Liu, X., Sternberg, T. and Zhai, L.: Satellite-based ground PM<sub>2.5</sub> estimation  
725 using timely structure adaptive modeling, *Remote Sens. Environ.*, 186, 152–163,  
726 doi:10.1016/j.rse.2016.08.027, 2016.





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

730 Fuzzi, S., Baltensperger, U., Carslaw, K., Decesari, S., Denier van der Gon, H., Facchini, M. C.,  
731 Fowler, D., Koren, I., Langford, B., Lohmann, U., Nemitz, E., Pandis, S., Riipinen, I., Rudich, Y.,  
732 Schaap, M., Slowik, J. G., Spracklen, D. V., Vignati, E., Wild, M., Williams, M. and Gilardoni, S.:  
733 Particulate matter, air quality and climate: lessons learned and future needs, *Atmos. Chem. Phys.*,  
734 15(14), 8217–8299, doi:10.5194/acp-15-8217-2015, 2015.

735 Gao, M., Beig, G., Song, S., Zhang, H., Hu, J., Ying, Q., Liang, F., Liu, Y., Wang, H., Lu, X.,  
736 Zhu, T., Carmichael, G. R., Nielsen, C. P. and McElroy, M. B.: The impact of power generation  
737 emissions on ambient PM<sub>2.5</sub> pollution and human health in China and India, *Environ. Int.*,  
738 121(August), 250–259, doi:10.1016/j.envint.2018.09.015, 2018.

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

743 Geng, G., Xiao, Q., Liu, S., Liu, X., Cheng, J., Zheng, Y., Xue, T., Tong, D., Zheng, B., Peng, Y.,  
744 Huang, X., He, K. and Zhang, Q.: Tracking Air Pollution in China: Near Real-Time PM<sub>2.5</sub> Retrievals  
745 from Multisource Data Fusion, *Environ. Sci. Technol.*, acs.est.1c01863, doi:10.1021/acs.est.1c01863,  
746 2021b.

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

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

754 He, Q., Gu, Y. and Zhang, M.: Spatiotemporal trends of PM<sub>2.5</sub> concentrations in central China  
755 from 2003 to 2018 based on MAIAC-derived high-resolution data, *Environ. Int.*, 137(August 2019),  
756 105536, doi:10.1016/j.envint.2020.105536, 2020.



- 757 Hu, X., Waller, L. A., Lyapustin, A., Wang, Y., Al-Hamdan, M. Z., Crosson, W. L., Estes, M. G.,  
758 Estes, S. M., Quattrochi, D. A., Puttaswamy, S. J. and Liu, Y.: Estimating ground-level PM<sub>2.5</sub>  
759 concentrations in the Southeastern United States using MAIAC AOD retrievals and a two-stage model,  
760 *Remote Sens. Environ.*, 140, 220–232, doi:10.1016/j.rse.2013.08.032, 2014.
- 761 Huang, C., Hu, J., Xue, T., Xu, H. and Wang, M.: High-Resolution Spatiotemporal Modeling for  
762 Ambient PM<sub>2.5</sub> Exposure Assessment in China from 2013 to 2019, *Environ. Sci. Technol.*, 55(3),  
763 2152–2162, doi:10.1021/acs.est.0c05815, 2021.
- 764 Jun, C., Ban, Y. and Li, S.: Open access to Earth land-cover map, *Nature*, 514(7523), 434–434,  
765 doi:10.1038/514434c, 2014.
- 766 Kolda, T. G. and Bader, B. W.: Tensor Decompositions and Applications, *SIAM Rev.*, 51(3), 455–  
767 500, doi:10.1137/07070111X, 2009.
- 768 de Leeuw, G., Sogacheva, L., Rodriguez, E., Kourtidis, K., Georgoulias, A. K., Alexandri, G.,  
769 Amiridis, V., Proestakis, E., Marinou, E., Xue, Y. and van der A, R.: Two decades of satellite  
770 observations of AOD over mainland China using ATSR-2, AATSR and MODIS/Terra: data set  
771 evaluation and large-scale patterns, *Atmos. Chem. Phys.*, 18(3), 1573–1592, doi:10.5194/acp-18-  
772 1573-2018, 2018.
- 773 Li, J., Li, C. and Zhao, C.: Different trends in extreme and median surface aerosol extinction  
774 coefficients over China inferred from quality-controlled visibility data, *Atmos. Chem. Phys.*, 18(5),  
775 3289–3298, doi:10.5194/acp-18-3289-2018, 2018a.
- 776 Li, L., Zhang, J., Meng, X., Fang, Y., Ge, Y., Wang, J., Wang, C., Wu, J. and Kan, H.: Estimation  
777 of PM<sub>2.5</sub> concentrations at a high spatiotemporal resolution using constrained mixed-effect bagging  
778 models with MAIAC aerosol optical depth, *Remote Sens. Environ.*, 217(January), 573–586,  
779 doi:10.1016/j.rse.2018.09.001, 2018b.
- 780 Li, L., Franklin, M., Girguis, M., Lurmann, F., Wu, J., Pavlovic, N., Breton, C., Gilliland, F. and  
781 Habre, R.: Spatiotemporal imputation of MAIAC AOD using deep learning with downscaling, *Remote*  
782 *Sens. Environ.*, 237(October 2019), 111584, doi:10.1016/j.rse.2019.111584, 2020.
- 783 Li, K., Bai, K., Li, Z., Guo, J. and Chang, N.-B.: Synergistic Data Fusion of Multimodal AOD and  
784 Air Quality Data for Near Real-Time Full Coverage Air Pollution Assessment, *J. Environ. Manage.*,  
785 2021. In revision
- 786 Li, Z., Zhang, Y., Shao, J., Li, B., Hong, J., Liu, D., Li, D., Wei, P., Li, W., Li, L., Zhang, F., Guo,



- 787 J., Deng, Q., Wang, B., Cui, C., Zhang, W., Wang, Z., Lv, Y., Xu, H., Chen, X., Li, L. and Qie, L.:  
788 Remote sensing of atmospheric particulate mass of dry PM<sub>2.5</sub> near the ground: Method validation  
789 using ground-based measurements, *Remote Sens. Environ.*, 173, 59–68, doi:10.1016/j.rse.2015.11.019,  
790 2016.
- 791 Li, Z., Wang, Y., Guo, J., Zhao, C., Cribb, M. C., Dong, X., Fan, J., Gong, D., Huang, J., Jiang,  
792 M., Jiang, Y., Lee, S. S., Li, H., Li, J., Liu, J., Qian, Y., Rosenfeld, D., Shan, S., Sun, Y., Wang, H.,  
793 Xin, J., Yan, X., Yang, X., Yang, X. qun, Zhang, F. and Zheng, Y.: East Asian Study of Tropospheric  
794 Aerosols and their Impact on Regional Clouds, Precipitation, and Climate (EAST-AIRCPC), *J.*  
795 *Geophys. Res. Atmos.*, 124(23), 13026–13054, doi:10.1029/2019JD030758, 2019.
- 796 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,  
797 X. and Yu, Q.: Estimation of long-term population exposure to PM<sub>2.5</sub> for dense urban areas using 1-  
798 km MODIS data, *Remote Sens. Environ.*, 179, 13–22, doi:10.1016/j.rse.2016.03.023, 2016.
- 799 Liu, M., Bi, J. and Ma, Z.: Visibility-Based PM<sub>2.5</sub> Concentrations in China: 1957–1964 and  
800 1973–2014, *Environ. Sci. Technol.*, 51(22), 13161–13169, doi:10.1021/acs.est.7b03468, 2017.
- 801 Liu, Y., Paciorek, C. J. and Koutrakis, P.: Estimating Regional Spatial and Temporal Variability  
802 of PM<sub>2.5</sub> Concentrations Using Satellite Data, Meteorology, and Land Use Information, *Environ.*  
803 *Health Perspect.*, 117(6), 886–892, doi:10.1289/ehp.0800123, 2009.
- 804 Lyapustin, A., Martonchik, J., Wang, Y., Laszlo, I. and Korokin, S.: Multiangle implementation of  
805 atmospheric correction (MAIAC): 1. Radiative transfer basis and look-up tables, *J. Geophys. Res.*  
806 *Atmos.*, 116(3), doi:10.1029/2010JD014985, 2011.
- 807 Lyapustin, A., Wang, Y., Korokin, S. and Huang, D.: MODIS Collection 6 MAIAC algorithm,  
808 *Atmos. Meas. Tech.*, 11(10), 5741–5765, doi:10.5194/amt-11-5741-2018, 2018.
- 809 Lyu, B., Hu, Y., Zhang, W., Du, Y., Luo, B., Sun, X., Sun, Z., Deng, Z., Wang, X., Liu, J., Wang,  
810 X. and Russell, A. G.: Fusion Method Combining Ground-Level Observations with Chemical  
811 Transport Model Predictions Using an Ensemble Deep Learning Framework: Application in China to  
812 Estimate Spatiotemporally-Resolved PM<sub>2.5</sub> Exposure Fields in 2014–2017, *Environ. Sci. Technol.*,  
813 53(13), 7306–7315, doi:10.1021/acs.est.9b01117, 2019.
- 814 Ma, Z., Hu, X., Huang, L., Bi, J. and Liu, Y.: Estimating Ground-Level PM<sub>2.5</sub> in China Using  
815 Satellite Remote Sensing, *Environ. Sci. Technol.*, 48(13), 7436–7444, doi:10.1021/es5009399, 2014.
- 816 Ma, Z., Hu, X., Sayer, A. M., Levy, R., Zhang, Q., Xue, Y., Tong, S., Bi, J., Huang, L. and Liu,



- 817 Y.: Satellite-based spatiotemporal trends in PM<sub>2.5</sub> concentrations: China, 2004-2013, *Environ. Health*  
818 *Perspect.*, 124(2), 184–192, doi:10.1289/ehp.1409481, 2016.
- 819 Park, S., Lee, J., Im, J., Song, C. K., Choi, M., Kim, J., Lee, S., Park, R., Kim, S. M., Yoon, J.,  
820 Lee, D. W. and Quackenbush, L. J.: Estimation of spatially continuous daytime particulate matter  
821 concentrations under all sky conditions through the synergistic use of satellite-based AOD and  
822 numerical models, *Sci. Total Environ.*, 713, 136516, doi:10.1016/j.scitotenv.2020.136516, 2020.
- 823 Shen, F., Zhang, L., Jiang, L., Tang, M., Gai, X., Chen, M. and Ge, X.: Temporal variations of six  
824 ambient criteria air pollutants from 2015 to 2018, their spatial distributions, health risks and  
825 relationships with socioeconomic factors during 2018 in China, *Environ. Int.*, 137(February), 105556,  
826 doi:10.1016/j.envint.2020.105556, 2020.
- 827 Sidiropoulos, N. D., De Lathauwer, L., Fu, X., Huang, K., Papalexakis, E. E. and Faloutsos, C.:  
828 Tensor Decomposition for Signal Processing and Machine Learning, *IEEE Trans. Signal Process.*,  
829 65(13), 3551–3582, doi:10.1109/TSP.2017.2690524, 2017.
- 830 Sogacheva, L., Popp, T., Sayer, A. M., Dubovik, O., Garay, M. J., Heckel, A., Christina Hsu, N.,  
831 Jethva, H., Kahn, R. A., Kolmonen, P., Kosmale, M., De Leeuw, G., Levy, R. C., Litvinov, P.,  
832 Lyapustin, A., North, P., Torres, O. and Arola, A.: Merging regional and global aerosol optical depth  
833 records from major available satellite products, *Atmos. Chem. Phys.*, 20(4), 2031–2056,  
834 doi:10.5194/acp-20-2031-2020, 2020.
- 835 Sun, J.-L., Jing, X., Chang, W.-J., Chen, Z.-X. and Zeng, H.: Cumulative health risk assessment  
836 of halogenated and parent polycyclic aromatic hydrocarbons associated with particulate matters in  
837 urban air, *Ecotoxicol. Environ. Saf.*, 113, 31–37, doi:10.1016/j.ecoenv.2014.11.024, 2015.
- 838 Sun, Z., Chang, N. Bin, Chen, C. F., Mostafiz, C. and Gao, W.: Ensemble learning via higher  
839 order singular value decomposition for integrating data and classifier fusion in water quality  
840 monitoring, *IEEE J. Sel. Top. Appl. Earth Obs. Remote Sens.*, 14, 3345–3360,  
841 doi:10.1109/JSTARS.2021.3055798, 2021.
- 842 Tang, Q., Bo, Y. and Zhu, Y.: Spatiotemporal fusion of multiple-satellite aerosol optical depth  
843 (AOD) products using Bayesian maximum entropy method, *J. Geophys. Res. Atmos.*, 121(8), 4034–  
844 4048, doi:10.1002/2015JD024571, 2016.
- 845 Tucker, L. R.: Some mathematical notes on three-mode factor analysis, *Psychometrika*, 31(3),  
846 279–311, doi:10.1007/BF02289464, 1966.



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

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

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

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

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

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

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

868 Wei, J., Li, Z., Lyapustin, A., Sun, L., Peng, Y., Xue, W., Su, T. and Cribb, M.: Reconstructing 1-  
869 km-resolution high-quality PM<sub>2.5</sub> data records from 2000 to 2018 in China: spatiotemporal variations  
870 and policy implications, *Remote Sens. Environ.*, 252(January 2020), 112136,  
871 doi:10.1016/j.rse.2020.112136, 2021a.

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

875 Wei, X., Bai, K., Chang, N. and Gao, W.: Multi-source hierarchical data fusion for high-resolution  
876 AOD mapping in a forest fire event, *Int. J. Appl. Earth Obs. Geoinf.*, 102(May), 102366,



- 877 doi:10.1016/j.jag.2021.102366, 2021b.
- 878 Xiao, Q., Zhang, H., Choi, M., Li, S., Kondragunta, S., Kim, J., Holben, B., Levy, R. C. and Liu,  
879 Y.: Evaluation of VIIRS, GOCI, and MODIS Collection 6 AOD retrievals against ground  
880 sunphotometer observations over East Asia, *Atmos. Chem. Phys.*, 16(3), 1255–1269, doi:10.5194/acp-  
881 16-1255-2016, 2016.
- 882 Xiao, Q., Wang, Y., Chang, H. H., Meng, X., Geng, G., Lyapustin, A. and Liu, Y.: Full-coverage  
883 high-resolution daily PM<sub>2.5</sub> estimation using MAIAC AOD in the Yangtze River Delta of China,  
884 *Remote Sens. Environ.*, 199(May), 437–446, doi:10.1016/j.rse.2017.07.023, 2017a.
- 885 Xiao, Q., Wang, Y., Chang, H. H., Meng, X., Geng, G., Lyapustin, A. and Liu, Y.: Full-coverage  
886 high-resolution daily PM<sub>2.5</sub> estimation using MAIAC AOD in the Yangtze River Delta of China,  
887 *Remote Sens. Environ.*, 199, 437–446, doi:10.1016/j.rse.2017.07.023, 2017b.
- 888 Xiao, Q., Chang, H. H., Geng, G. and Liu, Y.: An Ensemble Machine-Learning Model to Predict  
889 Historical PM<sub>2.5</sub> Concentrations in China from Satellite Data, *Environ. Sci. Technol.*,  
890 doi:10.1021/acs.est.8b02917, 2018.
- 891 Xiao, Q., Geng, G., Liang, F., Wang, X., Lv, Z., Lei, Y., Huang, X., Zhang, Q., Liu, Y. and He,  
892 K.: Changes in spatial patterns of PM<sub>2.5</sub> pollution in China 2000 – 2018 : Impact of clean air policies,  
893 *Environ. Int.*, 141(April), 105776, doi:10.1016/j.envint.2020.105776, 2020.
- 894 Xin, J., Wang, Y., Pan, Y., Ji, D., Liu, Z., Wen, T., Wang, Y., Li, X., Sun, Y., Sun, J., Wang, P.,  
895 Wang, G., Wang, X., Cong, Z., Song, T., Hu, B., Wang, L., Tang, G., Gao, W., Guo, Y., Miao, H.,  
896 Tian, S. and Wang, L.: The Campaign on Atmospheric Aerosol Research Network of China: CARE-  
897 China, *Bull. Am. Meteorol. Soc.*, 96(7), 1137–1155, doi:10.1175/BAMS-D-14-00039.1, 2015.
- 898 Xu, H., Guang, J., Xue, Y., de Leeuw, G., Che, Y. H., Guo, J., He, X. W. and Wang, T. K.: A  
899 consistent aerosol optical depth (AOD) dataset over mainland China by integration of several AOD  
900 products, *Atmos. Environ.*, 114, 48–56, doi:10.1016/j.atmosenv.2015.05.023, 2015.
- 901 Xue, T., Zheng, Y., Tong, D., Zheng, B., Li, X., Zhu, T. and Zhang, Q.: Spatiotemporal continuous  
902 estimates of PM<sub>2.5</sub> concentrations in China, 2000–2016: A machine learning method with inputs from  
903 satellites, chemical transport model, and ground observations, *Environ. Int.*, 123(December 2018),  
904 345–357, doi:10.1016/j.envint.2018.11.075, 2019.
- 905 Yang, F., Tan, J., Zhao, Q., Du, Z., He, K., Ma, Y., Duan, F., Chen, G. and Zhao, Q.:  
906 Characteristics of PM<sub>2.5</sub> speciation in representative megacities and across China, *Atmos. Chem.*



- 907 Phys., 11(11), 5207–5219, doi:10.5194/acp-11-5207-2011, 2011.
- 908 Yang, J. and Huang, X.: 30 m annual land cover and its dynamics in China from 1990 to 2019,  
909 Earth Syst. Sci. Data Discuss., 2021(April), 1–29, doi:<https://doi.org/10.5194/essd-2021-7>, 2021.
- 910 Yi-Lei Chen, Chiou-Ting Hsu and Liao, H.-Y. M.: Simultaneous Tensor Decomposition and  
911 Completion Using Factor Priors, IEEE Trans. Pattern Anal. Mach. Intell., 36(3), 577–591,  
912 doi:10.1109/TPAMI.2013.164, 2014.
- 913 Zhang, T., Zeng, C., Gong, W., Wang, L., Sun, K., Shen, H., Zhu, Z. and Zhu, Z.: Improving  
914 spatial coverage for Aqua MODIS AOD using NDVI-based multi-temporal regression analysis,  
915 Remote Sens., 9(4), doi:10.3390/rs9040340, 2017.
- 916 Zhang, Y., Gao, L., Cao, L., Yan, Z. and Wu, Y.: Decreasing atmospheric visibility associated  
917 with weakening winds from 1980 to 2017 over China, Atmos. Environ., 224(July 2019), 117314,  
918 doi:10.1016/j.atmosenv.2020.117314, 2020.
- 919 Zhao, C., Yang, Y., Fan, H., Huang, J., Fu, Y., Zhang, X., Kang, S., Cong, Z., Letu, H. and Menenti,  
920 M.: Aerosol characteristics and impacts on weather and climate over the Tibetan Plateau, Natl. Sci.  
921 Rev., 7(3), 492–495, doi:10.1093/nsr/nwz184, 2020.
- 922