

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

3 Kaixu Bai^{1,2*}, Ke Li¹, Mingliang Ma³, Kaitao Li⁴, Zhengqiang Li⁴, Jianping Guo^{5*},
4 Ni-Bin Chang⁶, Zhuo Tan¹, Di Han¹

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

7 ²Institute of Eco-Chongming, 20 Cuiniao Rd., Chongming, Shanghai 202162, China

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

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

11 ⁵State Key Laboratory of Severe Weather, Chinese Academy of Meteorological Sciences, Beijing, China

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

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

18 **Abstract.** Developing a big data analytics framework for generating a Long-term Gap-free High-
19 resolution Air Pollutants concentration dataset (abbreviated as LGHAP) is of great significance for
20 environmental management and earth system science analysis. By synergistically integrating
21 multimodal aerosol data acquired from diverse sources via a tensor flow based data fusion method, a
22 gap-free aerosol optical depth (AOD) dataset with daily 1-km resolution covering the period of 2000–
23 2020 in China was generated. Specifically, data gaps in daily AOD imageries from MODIS aboard
24 Terra were reconstructed based on a set of AOD data tensors acquired from diverse satellites,
25 numerical analysis, and *in situ* air quality measurements via integrative efforts of spatial pattern
26 recognition for high dimensional gridded image analysis and knowledge transfer in statistical data
27 mining. To our knowledge, this is the first long-term gap-free high resolution AOD dataset in China,
28 from which spatially contiguous PM_{2.5} and PM₁₀ concentrations were then estimated using an
29 ensemble learning approach. Ground validation results indicate that the LGHAP AOD data are in a
30 good agreement with *in situ* AOD observations from AERONET, with R of 0.91 and RMSE equaling
31 to 0.21. Meanwhile, PM_{2.5} and PM₁₀ estimations also agreed well with ground measurements, with R
32 of 0.95 and 0.94 and RMSE of 12.03 and 19.56 $\mu\text{g m}^{-3}$, respectively. The LGHAP provides a suite of
33 long-term gap free gridded maps with high-resolution to better examine aerosol changes in China over
34 the past two decades, from which three major variation periods of haze pollution were revealed in
35 China. Additionally, the proportion of population exposed to unhealthy PM_{2.5} was increased from
36 50.60% in 2000 to 63.81% in 2014 across China, which was then reduced drastically to 34.03% in
37 2020. Overall, the generated LGHAP dataset has a great potential to trigger multidisciplinary
38 applications in earth observations, climate change, public health, ecosystem assessment, and
39 environmental management. The daily resolution AOD, PM_{2.5}, and PM₁₀ datasets are publicly
40 available at <https://doi.org/10.5281/zenodo.5652257> (Bai et al., 2021a),
41 <https://doi.org/10.5281/zenodo.5652265> (Bai et al., 2021b), and
42 <https://doi.org/10.5281/zenodo.5652263> (Bai et al., 2021c), respectively. Monthly and annual datasets
43 can be acquired from <https://doi.org/10.5281/zenodo.5655797> (Bai et al., 2021d) and
44 <https://doi.org/10.5281/zenodo.5655807> (Bai et al., 2021e), respectively. Python, Matlab, R, and IDL
45 codes were also provided to help users read and visualize these data.

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

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; Yang et al., 2020; Zheng et al., 2020). Monitoring aerosol loading in the
52 atmosphere is thus of great significance for climate change attribution and haze pollution assessment.
53 Aerosol optical depth (AOD), an indicator of aerosol bulks distributed within a column of air from the
54 Earth's surface to the top of the atmosphere, has been monitored for decades to map global aerosol
55 loading in the atmosphere. Compared with sparsely and unevenly distributed ground-based aerosol
56 monitoring stations (e.g., AERONET), satellite instruments can map AOD with vaster spatial coverage
57 at even sub-hourly sampling frequency (e.g., geostationary satellite). An overview of sensors,
58 algorithms, and AOD datasets that are widely used in the community can be found in the literature
59 such as Sogacheva et al. (2020) and Wei et al. (2020).

60 Due to negative impacts of bright surface (e.g., snow cover) and clouds, as well as algorithmic
61 restrictions, satellite AOD retrievals often suffer from extensive data gaps, significantly reducing the
62 downstream application potential such as mapping particulate matter (PM) concentrations at the
63 ground surface (e.g., Bai et al., 2019a; Wei et al., 2021a). Also, excessive data gaps in AOD imageries
64 may result in large uncertainty when assessing aerosol impacts on weather and climate (Guo et al.,
65 2017; Li et al., 2019; Zhao et al., 2020; Zheng et al., 2018). Over the years, versatile gap filling methods
66 have been developed (e.g., Bai et al., 2016, 2020b; Chang et al., 2015). Nonetheless, filling data gaps
67 in satellite-based AOD retrievals is still challenging due to extraordinary nonrandom missing values
68 and high aerosol dynamics in space and time.

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

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

80 Although there exist a variety of gap filling methods, spatially gap free AOD datasets are still
81 rare, particularly high-resolution AOD datasets from satellites, significantly limiting downstream
82 applications such as PM_x concentration mapping. In spite of versatile $PM_{2.5}$ concentration prediction
83 models (e.g., Di et al., 2019; Fang et al., 2016; Hu et al., 2014; Li et al., 2016; Lin et al., 2016; Liu et
84 al., 2009; Wang et al., 2021a), to date, there are few publicly accessible PM_x concentration datasets
85 that can be used to examine haze pollution variations regionally and globally. Several typical datasets,
86 e.g., the one generated by the Dalhousie University (van Donkelaar et al., 2010, 2016), CHAP (Wei et
87 al., 2021a), and TAP (Geng et al., 2021), have been widely applied to advance our understanding on
88 aerosol impacts across China and globe. However, these datasets more or less still suffer from
89 drawbacks in spatial and/or temporal resolution, spatial coverage, and data accuracy. To meet the
90 contemporary needs, Zhang et al. (2021) provided a more comprehensive review of the widely used
91 PM_x concentration mapping approaches. With a thorough review for $PM_{2.5}$ concentration mapping
92 techniques, an optimal full-coverage $PM_{2.5}$ modeling scheme was proposed, in which diverse aerosol
93 datasets were fused toward a full-coverage AOD map based on a multi-modal approach (Bai et al.,
94 2022). In parallel with these efforts, some attempted to improve AOD data coverage over space with
95 high accuracy by merging AODs observed at adjacent times directly (Li et al., 2022).

96 With such prior knowledge, the current study developed a big data analytics framework for
97 generating a Long-term Gap-free High-resolution Air Pollutants concentration dataset (abbreviated as
98 LGHAP hereafter), aiming at providing gap-free AOD, $PM_{2.5}$ and PM_{10} concentration data with a daily
99 1-km resolution in China for the period of 2000 to 2020. Toward such a goal, multimodal aerosol data
100 acquired from diverse sources including satellites, ground stations and numerical models were
101 synergistically integrated via the higher order singular value decomposition (HOSVD) to form a tensor
102 flow based data fusion framework in the current study. Full coverage $PM_{2.5}$ and PM_{10} concentration
103 data were then estimated on the basis of the gap-filled AOD dataset. This 21-year-long gap-free high
104 resolution (daily/1km) aerosol dataset was then compared against ground-based AOD and PM_x
105 observations to validate the data accuracy of each product, particularly their performance in spatial
106 pattern recognition and temporal trend assessment. These advances endorsed a better assessment of

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

109 2 Data sources

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

115 **Table 1.** Summary of the data sources used in this study to generate gap free high resolution AOD
 116 and PM_x 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 _{2.5} , PM ₁₀ , SO ₂ , NO ₂	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°

117 **2.1 Gridded aerosol products**

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

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

144 In addition to satellite-based AOD products, numerically simulated aerosol diagnostics from
145 MERRA-2, including AOD and aerosol components such as black carbon, organic carbon, dust and
146 sulfate, were also applied to help reconstruct missing AOD information and to predict PM_{2.5} and PM₁₀

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

150 **2.2 *In situ* AOD and air quality measurements**

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

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

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

176 visibility data was examined using an outlier detection method, i.e., the annual mean should not exceed
177 3 times the standard deviation of data over a 5-year time window (Zhang et al., 2020). Those with
178 apparent jumps and drifts in visibility time series were excluded. Meanwhile, visibility data on
179 rainstorm and foggy days were eliminated as well.

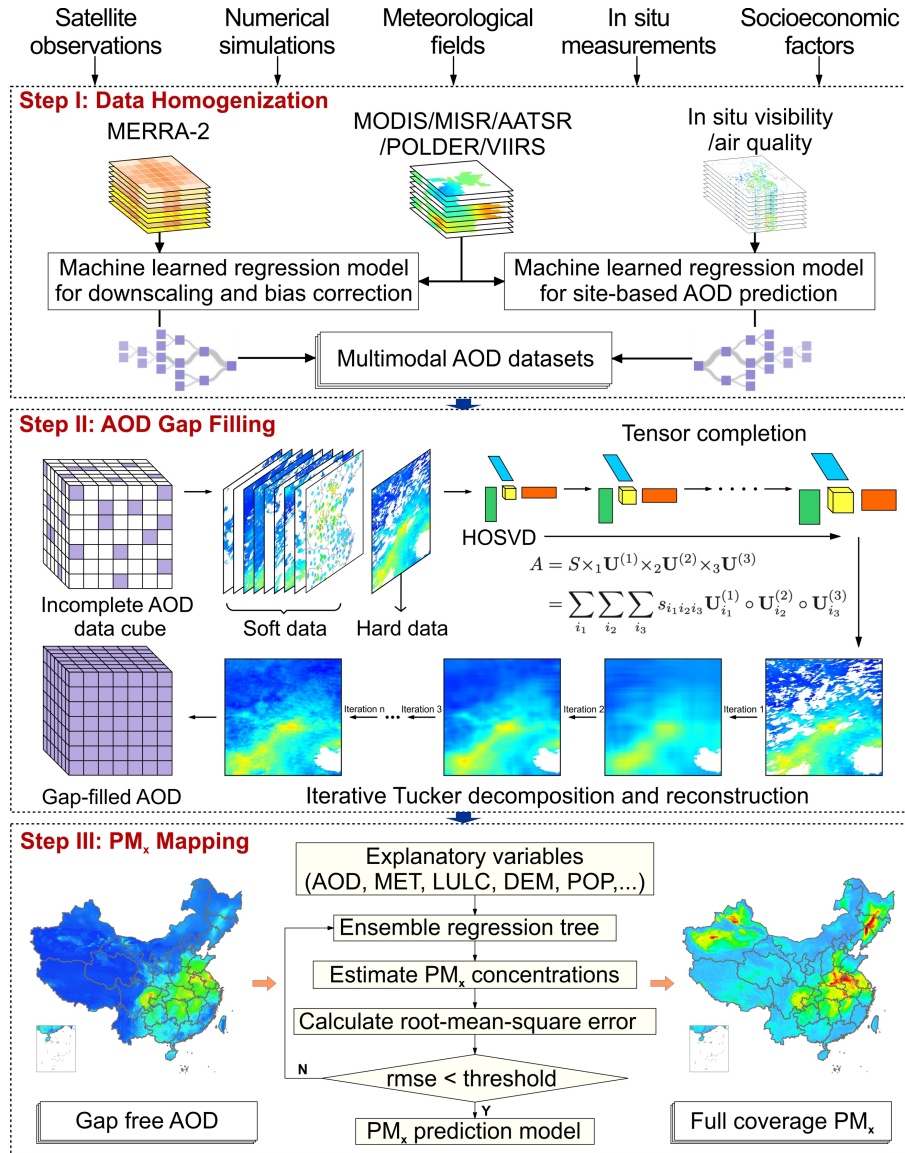
180 **2.3 Auxiliary data**

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

194 **3 Methodology**

195 Toward the generation of LGHAP aerosol datasets to advance environment management and
196 earth system science analysis, here we developed a big data analytics framework via a seamless
197 integration of the tensor flow based multimodal data fusion with ensemble learning based PM_x
198 concentration estimation. The proposed method transformed a set of data tensors of AOD and other
199 related datasets such as air pollutants concentration and atmospheric visibility that were acquired from
200 diversified sensors or platforms via integrative efforts of spatial pattern recognition for high
201 dimensional gridded data analysis toward data fusion and multiresolution image analysis, as well as
202 knowledge transfer in statistical data mining. The proposed method consists of three major procedures
203 in general, including multisensory data homogenization, tensor flow based AOD reconstruction, and

204 ensemble learning for PM_x concentration estimation. The analytical framework of the big data
 205 analytics is depicted in Figure 1 and described in details in the following subsections.



206
 207 **Figure 1.** Flowchart of the proposed big data analytics framework for generating a long-term gap-free
 208 high-resolution air pollutants concentration dataset (LGHAP), taking aerosol optical depth (AOD) and
 209 PM_x ($PM_{2.5}$ and PM_{10}) concentration in China as illustration. HOSVD is an acronym of high order
 210 singular value decomposition. MET, LULC, DEM, and POP denote variables of meteorology, land
 211 use/land cover, digit elevation model, and population, respectively.

212 3.1 Multisensory data homogenization

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

215 differences between these diversified datasets is crucial to multisensory data integration. In this study,
216 machine-learned regression models were established to harmonize these heterogeneous aerosol
217 datasets. A baseline dataset was first selected to be used as the learning target while other datasets
218 were calibrated to the level of baseline dataset to make them comparable. Given finer resolution and
219 higher proportion of data coverage in space and time, the MAIAC AOD product from Terra (AOD_{Terra})
220 was selected as the baseline dataset. Consequently, six machine-learned regression models were
221 established between AOD_{Terra} and each gridded AOD product (i.e., five satellite-based AOD products
222 plus MERRA-2 AOD simulations) using the random forest method. Meteorological factors (MET),
223 land cover types (LULC), topographic (DEM) and population (POP) were used as covariates to help
224 downscale these multimodal AOD products to have a resolution same as AOD_{Terra} while accounting
225 for cross-mission biases arising from temporal and algorithmic differences.

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

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

245 3.2 Tensor flow based AOD reconstruction

246 The core of generating full coverage AOD imageries is to fill in data gaps in AOD_{Terra} . Previous
247 studies have demonstrated that merging satellite AOD retrievals at adjacent time steps can help
248 improve the observational AOD coverage at each single snapshot, while the involvement of numerical
249 AOD simulations can help bridge AOD data gaps (Li et al., 2022; Bai et al., 2022). In this study, a
250 tensor completion method was particularly designed and applied to fulfil the gap filling in AOD_{Terra} .
251 Specifically, the incomplete AOD_{Terra} imageries were deemed as the hard data (true AOD state) while
252 other AOD datasets (e.g., the downscaled AOD datasets and site-specific AOD predictions inferred
253 from air pollutants concentration and atmospheric visibility) were used as the soft data (complementary
254 data) to help reconstruct AOD distribution in AOD_{Terra} via tensor flow based pattern recognition.
255 Detailed procedures for gap filling are outlined as follows.

256 3.2.1 Initial AOD tensor construction

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

267 For the screening of historical observations resembling AOD_{Terra} distribution on the given date
268 to be reconstructed, AOD_{M2} was used in concert with AOD_{Terra} and site-based AOD estimations to
269 identify similar imageries. Toward this goal, site-specific AOD estimations and 5% randomly selected
270 downscaled AOD_{M2} data were merged directly with valid AOD_{Terra} to form a new image on each date.
271 Subsequently, correlations and biases were estimated between AOD_{Terra} on the given date to be
272 reconstructed and each newly merged historical AOD_{Terra} image. To avoid the inclusion of imageries
273 with distinct variation patterns, only those closely resembling AOD_{Terra} on the date to be reconstructed

274 were finally retained in terms of their correlations and biases subject to a threshold of $R > 0.7$ and
275 $RMSE < 0.2$. Once sufficient historical imageries were obtained, the data tensor of AOD was
276 constructed by compiling the observed AOD imageries on the given date with historical imageries to
277 a three-dimension data array $\mathbf{A} \in \mathbf{R}^{N_1 \times N_2 \times N_3}$ (composed of N_3 images with a size of $N_1 \times N_2$).
278 Considering satellite AOD retrievals suffer from extensive data gaps, we injected data values of site-
279 specific AOD estimations and 1% randomly selected downscaled AOD_{M2} data directly onto grids
280 where AOD_{Terra} values missed on each specific date as prior knowledge. This not only accelerates
281 convergence speed during the reconstruction process but avoids large reconstruction errors over
282 regions with tremendous data gaps in satellite observed AOD imageries.

283 3.2.2 Gap filling via tensor completion

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

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

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

303 a.k.a, ε_1 in Table 2). Meanwhile, to make the computational burden manageable, the study region
 304 (China in this study) was divided into 40 subregions (refer to Figure S2 for the spatial distribution of
 305 these subregions), and the tensor completion was then performed over each individual region. Finally,
 306 the reconstructed imageries were mosaiced to attain a national gap-free AOD map on each specific
 307 date. During this step, a smooth filter was applied to solve the boundary effect when mosaicking two
 308 adjacent maps. Specifically, data value on each overlapped grid at the boundary (50 km on the edge of
 309 subregion) was averaged via an inverse distance (the distance to the edge) weighting scheme. In the
 310 end, the mosaic AOD_{Terra} image was retained as the final gap-free AOD product.

311 **Table 2.** The proposed tensor completion algorithm for AOD distribution reconstruction in AOD_{Terra}.

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

312 3.3 PM_x concentration estimation

313 In this study, the widely used random forest method was applied to establish regression models
 314 for PM_{2.5} and PM₁₀ concentration estimation. Ground measured PM_{2.5} (or PM₁₀) concentration data
 315 were used as the learning target while gap filled AOD, aerosol components (AER_{comp}), meteorological
 316 factors (MET), digital elevation model (DEM), NDVI, land cover information (LC), and population
 317 were used as regressors. The random forest regression model can be generally formulated as:

$$318 \quad \text{PM}_x = \text{RF}(\text{AOD}, \text{AER}_{\text{comp}}, \text{MET}, \text{DEM}, \text{NDVI}, \text{POP}, \text{LC}, \text{month}) \quad (1)$$

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

328 **3.4 Point-surface data fusion**

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

343 **4 Results and discussion**

344 **4.1 Data accuracy of gap-free AOD in LGHAP**

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

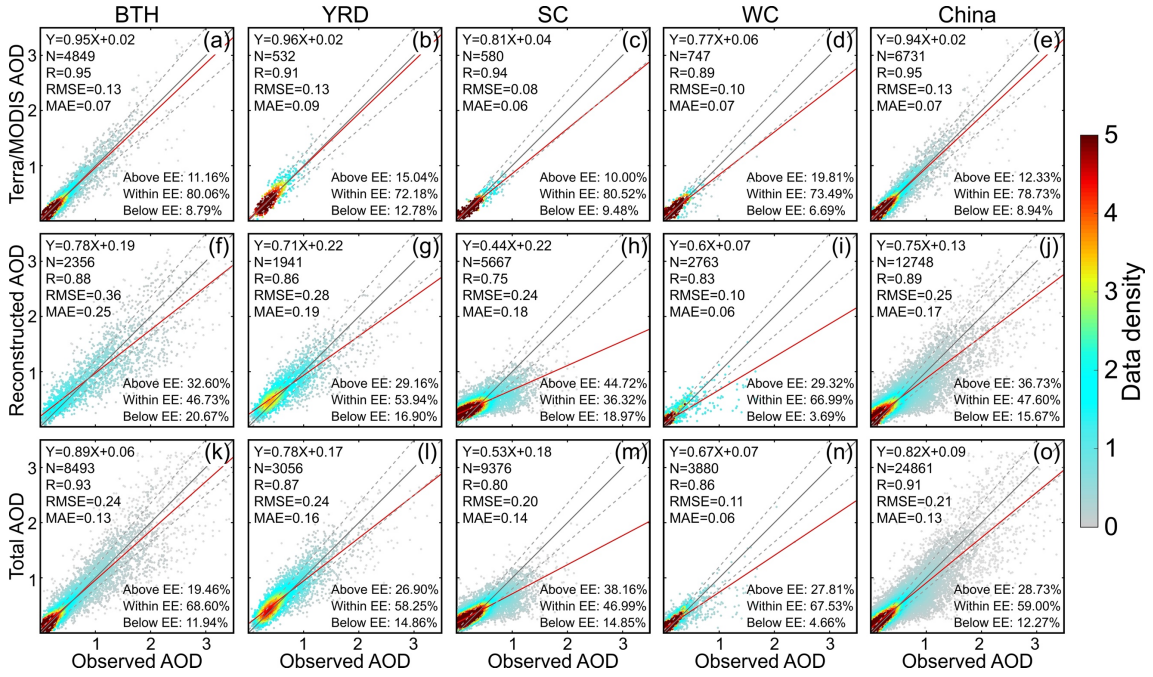
347 measurements were not used as data input when reconstructing missing AOD information, thereby the
 348 gap-free AOD can be directly compared with *in situ* AOD measurements from AERONET. As
 349 indicated, all these AOD datasets are in a good agreement with *in situ* AOD measurements. Generally,
 350 AODs from MODIS onboard Terra and Aqua have an almost identical data accuracy, which is also
 351 among the highest when comparing with other datasets (R=0.95 and RMSE=0.14). AODs from
 352 AATSR show a comparable accuracy with that of MODIS, but with a relatively low correlation with
 353 ground-based AOD measurements. AODs from MISR, POLDER and VIIRS exhibit a similar bias
 354 level, with R varying from 0.80 to 0.92 and RMSE ranging from 0.22 to 0.29. In contrast, AOD_{M2} data
 355 have the poorest accuracy among these eight gridded AOD datasets (R=0.77 and RMSE=0.36), even
 356 though AOD data from AERONET and satellite observations like MODIS had been already
 357 assimilated. This indicates the presence of large biases in AOD_{M2} and thus these AOD_{M2} data cannot
 358 be used solely to delineate AOD distributions over space.

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

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

361
 362 Compared to the first seven gridded AOD datasets, the LGHAP AOD dataset has an accuracy
 363 slightly worse than the original MODIS AOD product but comparable to AODs from MISR, POLDER
 364 and MERRA-2, with R of 0.91 and RMSE equaling to 0.21 compared to ground-based AOD
 365 observations. Nevertheless, the gap-filled AOD appeared to overestimate ground-based AOD
 366 observations, and this could be due to the involvement of AODs from VIIRS and POLDER as these
 367 two products significantly overestimated ground AOD observations, which can be indicated by the

368 proportion of data pairs above the expected error (EE). On the other hand, significant underestimations
369 in AOD_{M2} were not introduced to the LGHAP AOD as the former had a below EE ratio of 32.97%
370 which was only 12.27% in the latter. These results indicate that the LGHAP AOD data are more likely
371 to resemble AOD distributions revealed by satellite observations rather than AOD_{M2}, endorsing the
372 advantages of involving multisensory satellite AOD observations to support missing AOD
373 reconstruction. Figure 2 further compares the data accuracy of original AOD_{Terra} and the reconstructed
374 data over different regions of China. It is indicative that the purely reconstructed data have an accuracy
375 ($R=0.88$ and $RMSE=0.26$) lower than the original AOD_{Terra} ($R=0.95$ and $RMSE=0.13$) across China,
376 especially in South China where the reconstructed data were significantly underestimated the ground-
377 based AOD observations. Possible reasons for this effect could be attributed to extensive data gaps in
378 satellite AOD retrievals due to frequent and extensive cloud covers over there (refer to Figure S3 for
379 the distribution of mean data integrity of AOD_{Terra} during 2000–2020), and the scarce AOD
380 observations significantly limit the learning capacity in space and temporal domain during the tensor
381 completion process. In other words, limited observations in satellite imageries greatly reduced the
382 learning performance from the sparse tensor. Even though, the purely reconstructed data exhibit a bias
383 level comparable to AOD retrievals from several satellite instruments, e.g., MISR, VIIRS, and
384 POLDER. This demonstrates the good performance of the proposed tensor completion method in
385 reconstructing missing AOD information. By combining the reconstructed data with original AOD_{Terra},
386 we obtained a 21-year-long gap free high-resolution (daily/1-km) AOD dataset with satisfying
387 accuracy ($R=0.91$ and $RMSE=0.21$).



388

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

393

394

395

396

397

398

399

400

401

402

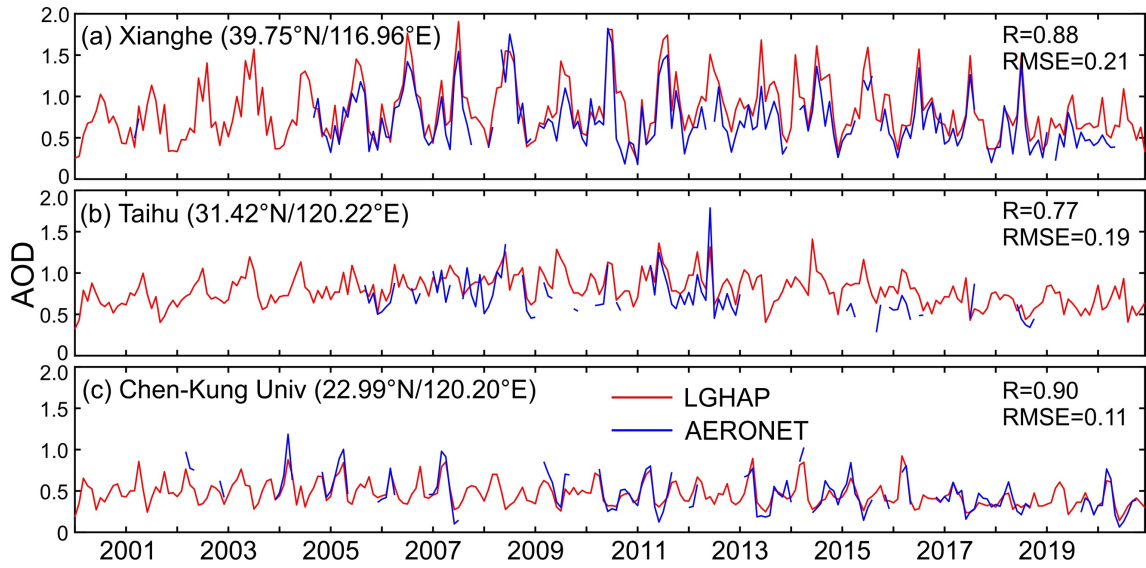
403

404

405

406

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

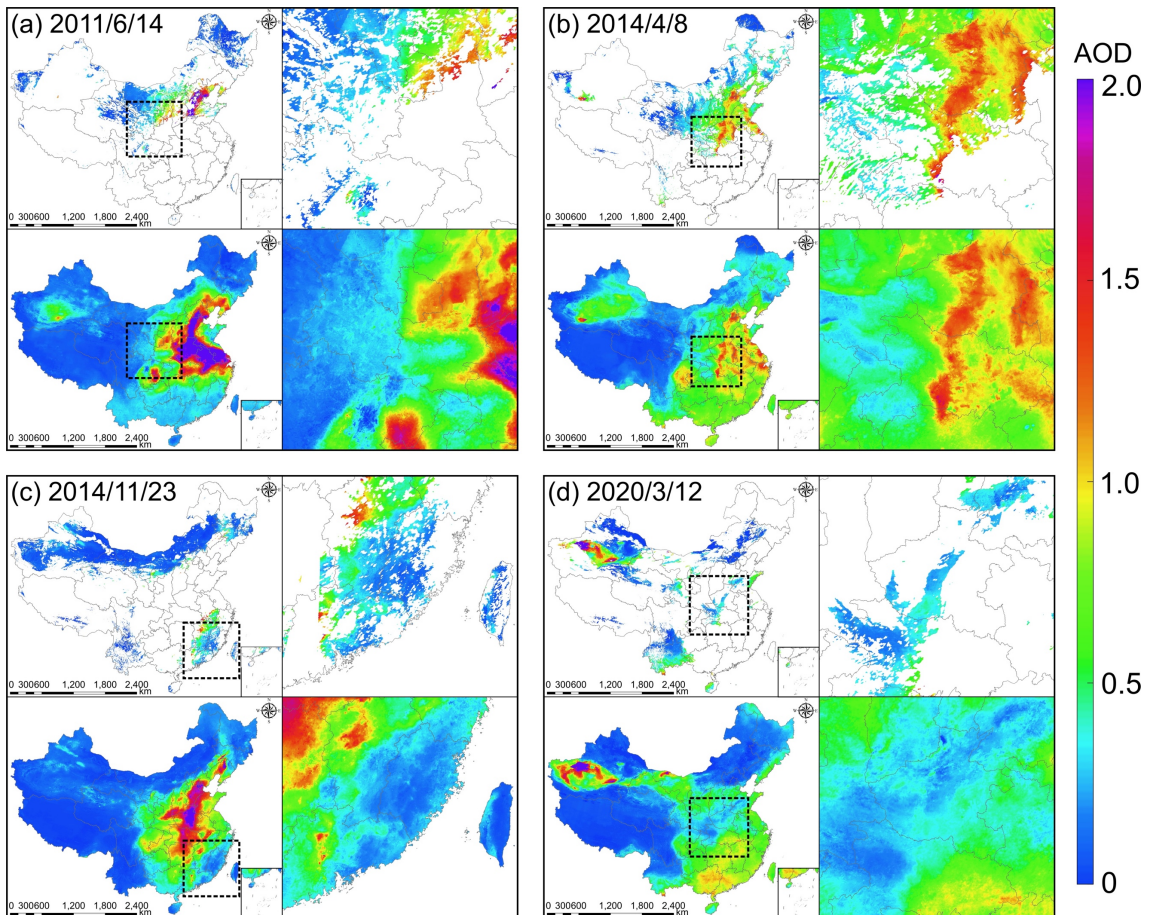


407

408

409

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



410

411

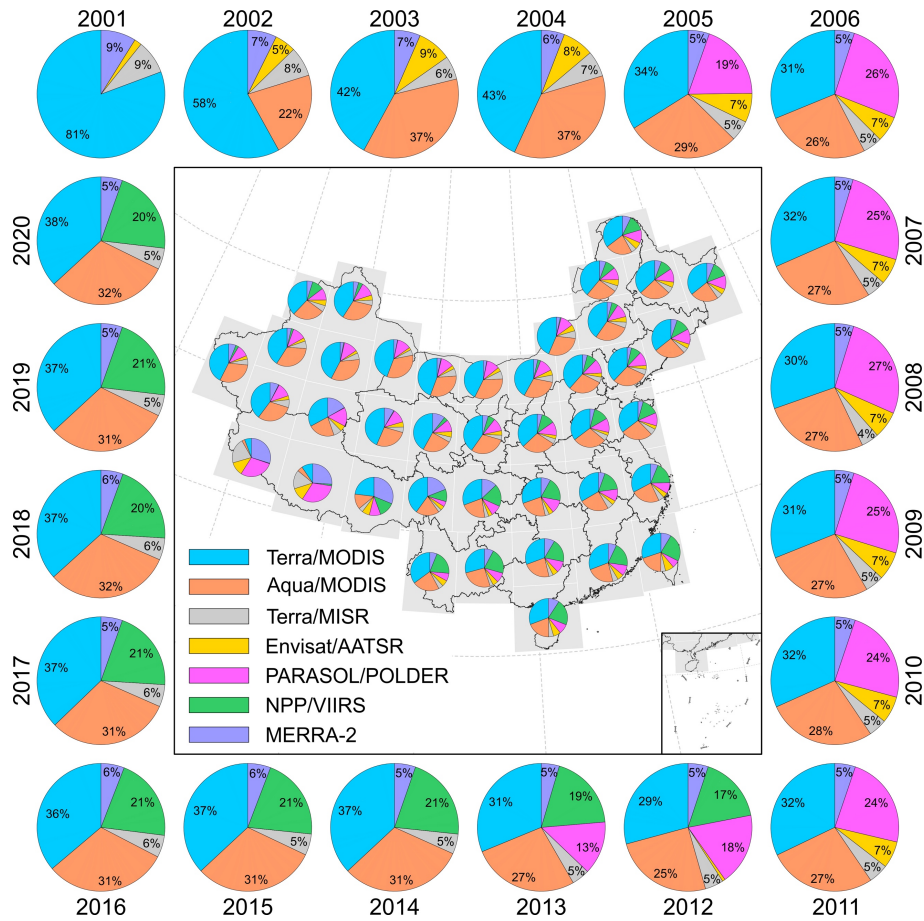
412

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

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

415 Since the final gap-free AOD product was generated mainly by integrating a set of data tensor
416 of gridded AOD with AOD estimations from *in situ* air quality measurements, the relative contribution
417 of each product to the final gap-free dataset is worth being investigated. In this study, a data coverage
418 ratio weighted nonlinear correlation coefficient was proposed to examine the relative contribution of
419 each gridded product to the LGHAP AOD dataset. The nonlinear correlation coefficient was used to
420 assess the mutual information between two variables (Sun et al., 2021; Wang et al., 2005), while the
421 data coverage ratio was multiplied to indicate the overall contribution of one product to the final fused
422 dataset (refer to Text S2 for the definition of this indicator). As shown in Figure 5, the relative
423 contribution of each gridded product varied with time and the input data sources. In the early two years
424 (2000–2001), the AOD distribution in gap-free imageries was determined largely by AOD_{Terra} (81%),
425 whereas this ratio decreased to about 30% when many other products were involved, especially AOD
426 from Aqua and PARASOL. With the advent of VIIRS and the loss of PARASOL after 2012, the
427 relative contribution changed drastically as AOD from MODIS and VIIRS played the dominant roles
428 in reconstructing AOD distribution. Note the relative contribution of AOD_{M2} remained lower than 10%,
429 indicative of the greater importance of satellite observations in generating the LGHAP AOD product.

430 With respect to the temporally averaged contribution in each subregion, it shows that the
431 relative contribution of each product also varied significantly across regions. Generally, AOD from
432 MODIS aboard Terra and Aqua played the most important role (>60%) in generating the LGHAP
433 AOD product, except over the southwest part of the country (Tibet plateau) where AOD_{M2} contributed
434 most. This is largely associated with the fact that data gaps are abnormally high in satellite observations
435 over this region because of the vast and long-lasting snow cover (refer to Figure S3 for the data
436 integrity distribution). Consequently, AOD_{M2} would play an important role in reconstructing AOD
437 distribution over such regions. Note that the relative contribution of AOD estimations from *in situ* air
438 quality measurements were not accounted for in the current analysis because of incomparable spatial
439 coverage of *in situ* data contrast to gridded AOD products, and this does not imply the contribution of
440 *in situ* AOD estimations being negligible. Overall, the results shown here clearly highlight the success
441 of big data analytics in generating the LGHAP AOD dataset via integrative efforts from diversified
442 data sources.



443

444

445

446

447

448

449

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

450

4.2 Data accuracy of PM_{2.5} and PM₁₀ estimations

451

452

453

454

455

456

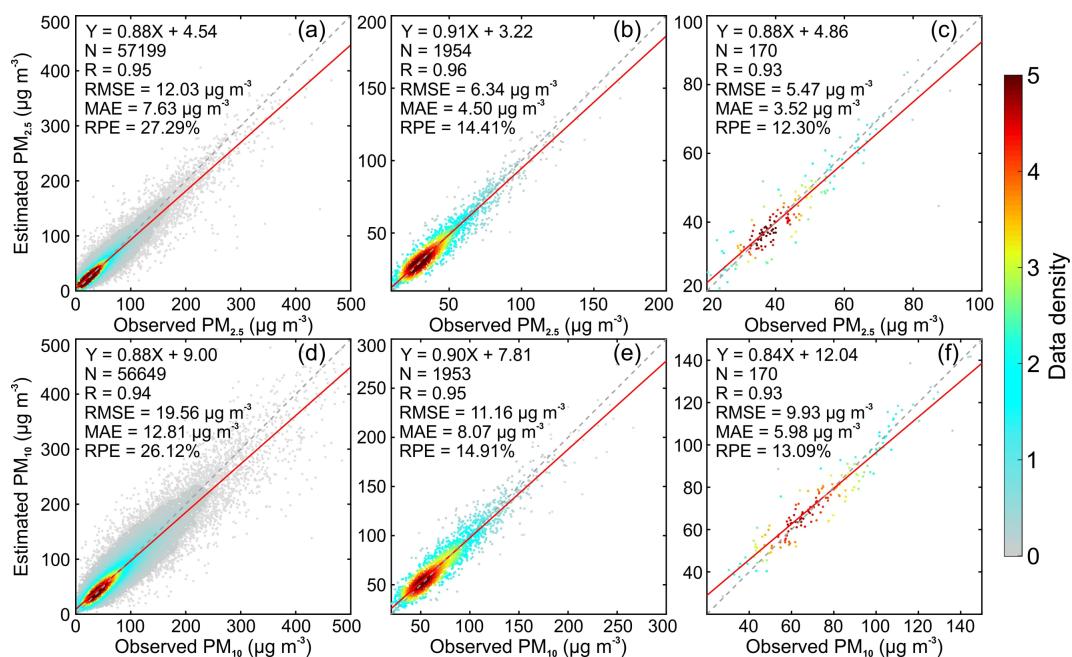
457

By taking advantage of the gap-filled AOD, daily 1-km resolution PM_{2.5} and PM₁₀ concentration data in China were then estimated via an ensemble learning approach. Figure S4 shows the sample-based cross validation accuracy of two prediction models. It shows that the original daily PM_{2.5} prediction model had a sample-based cross validation R² of 0.79 and RMSE of 20.04 μg m⁻³. This accuracy is comparable to our previous study (Bai et al., 2019a), but slightly worse than those reported in some recent studies (Table 4). In contrast, PM₁₀ had a much higher prediction accuracy, with R² of 0.90 and RMSE of 21.06 μg m⁻³ for the daily product. This good performance should be attributed to

458 the involvement of PM_{2.5} estimations as a predictor in the PM₁₀ prediction model. Figure 6 shows the
 459 site-specific (held-out in advance) validation accuracy of daily, monthly, and annual mean PM_{2.5} and
 460 PM₁₀ concentration in LGHAP. As shown, the site-specific validation results indicated that the final
 461 full-coverage (gap free) daily PM_{2.5} and PM₁₀ concentration data are in a good agreement with ground-
 462 based measurements, with R of 0.95 and RMSE of 12.03 $\mu\text{g m}^{-3}$ for PM_{2.5} while R of 0.94 and RMSE
 463 of 19.56 $\mu\text{g m}^{-3}$ for PM₁₀. Overall, PM_x data in LGHAP are not only spatially complete with a finer
 464 resolution but have a comparable accuracy with previous studies.

465 **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. (2021a)	No	1 km	2000~2018	0.86~0.90	10.09~18.39
Geng et al. (2021)	Yes	10 km	2000~2021	0.80~0.88	13.90~22.10
Xue et al. (2019)	Yes	10 km	2000~2016	0.61	27.80
Chen et al. (2018)	No	10 km	2005~2016	0.83	28.10
Lyu et al. (2019)	Yes	12 km	2014~2017	0.64	24.80
Ma et al. (2016)	No	10 km	2004~2013	0.79	27.42
Huang et al. (2021)	No	1 km	2013~2019	0.88	15.73
Xiao et al. (2018)	Yes	10 km	2013~2017	0.79	21.00
LGHAP PM _{2.5}	Yes	1 km	2000~2020	0.90	12.03

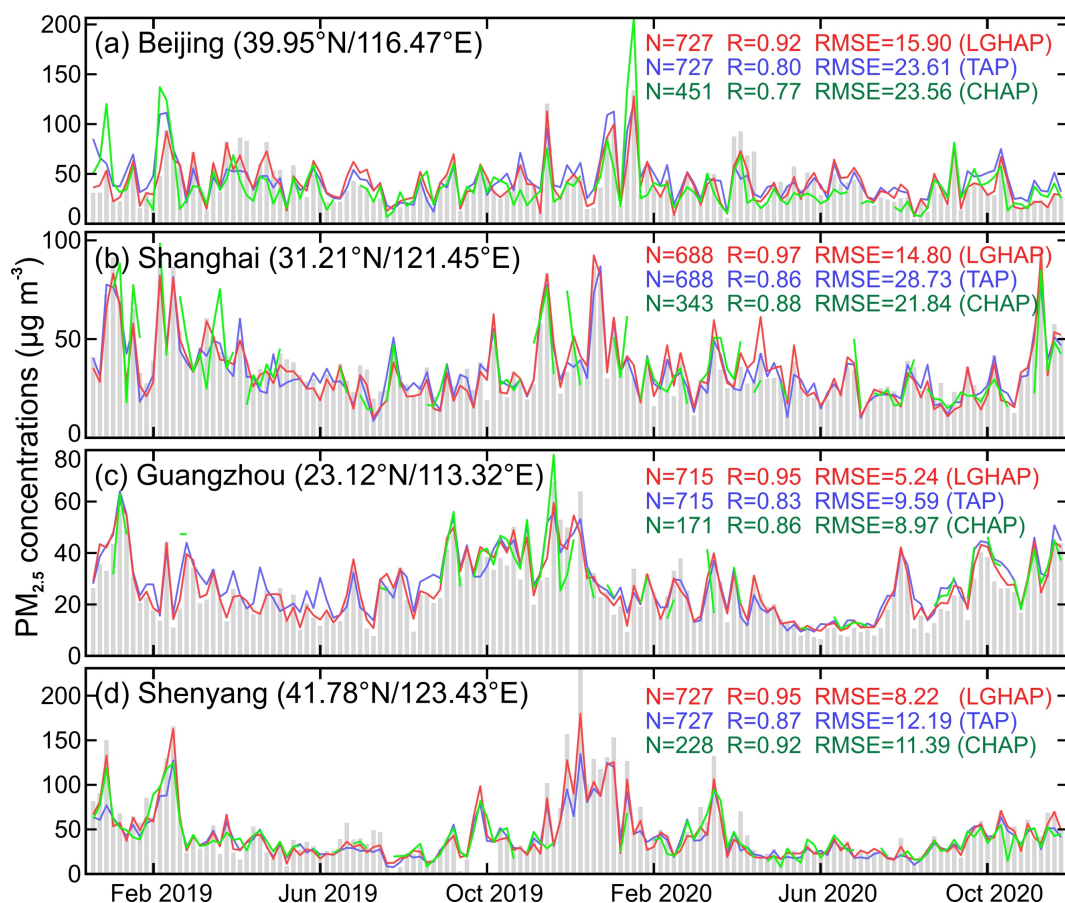


466

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

471

472 Figure 7 presents a two-year-long comparison of $PM_{2.5}$ concentration time series from LGHAP
473 and two other open access datasets with $PM_{2.5}$ measurements sampled at four United States Embassy
474 in China. Since this ground-based dataset has been seldomly noticed and used, it can be applied as an
475 independent dataset to fairly evaluate the accuracy of these three machine-learned $PM_{2.5}$ estimations.
476 As shown, all these three datasets well reconstructed temporal variations of $PM_{2.5}$ from 2019 to 2020.
477 Temporally, LGHAP and TAP are continuous while CHAP suffers from significant data gaps because
478 no gap filling was applied when generating the dataset. Compared with the other two datasets, LGHAP
479 $PM_{2.5}$ data had a better agreement with ground-based $PM_{2.5}$ measurements. This high accuracy could
480 be partially due to the fusion of *in situ* $PM_{2.5}$ data measured at adjacent sites via the OI method. Figure
481 S5 compares $PM_{2.5}$ time series from LGHAP with $PM_{2.5}$ measurements sampled at five United States
482 Embassy in China. It is indicative that historical $PM_{2.5}$ variations over these five cities were well
483 reconstructed in LGHAP, even over years before 2014 at which $PM_{2.5}$ measurements from state-
484 control monitoring sites were not available. Note $PM_{2.5}$ estimations appeared to significantly
485 underestimate $PM_{2.5}$ concentration sampled at the Embassy in Beijing before 2013. Considering the
486 reconstructed AOD time series agreed well with AERONET AOD in Beijing (Figure 3a), and the
487 model performed well in predicting historical $PM_{2.5}$ in Shanghai during the synchronous time period
488 (Figure S5b), we are more willing to attribute this issue to significant $PM_{2.5}$ overestimations by the US
489 Embassy during that period. Overall, these independent validation results collectively indicate a good
490 accuracy of $PM_{2.5}$ in LGHAP dataset.



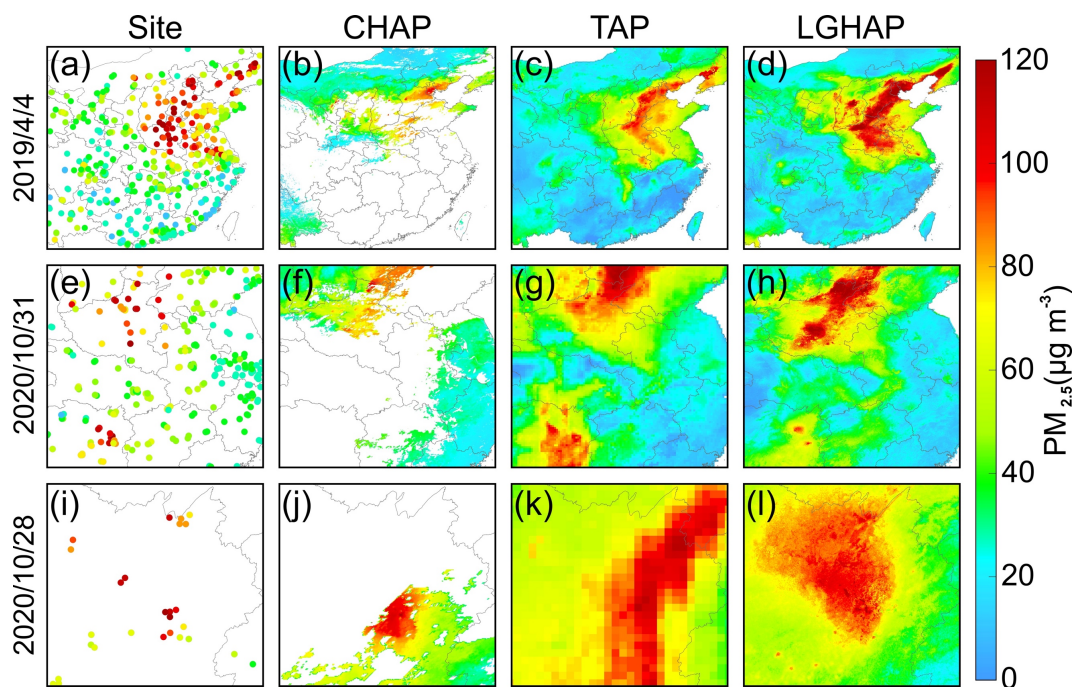
491

492 **Figure 7.** Comparison of PM_{2.5} concentration time series between LGHAP (red line) and two open
 493 datasets (blue: TAP, green: CHAP). Here, hourly PM_{2.5} concentrations measured by four United States
 494 Embassy in China from 2019 to 2020 (grey bar) were used as an independent PM_{2.5} dataset to validate
 495 these three daily products. CHAP and TAP are two open access datasets providing PM_{2.5}
 496 concentration that were created by Wei et al. (2021a) and Geng et al. (2021) respectively.

497

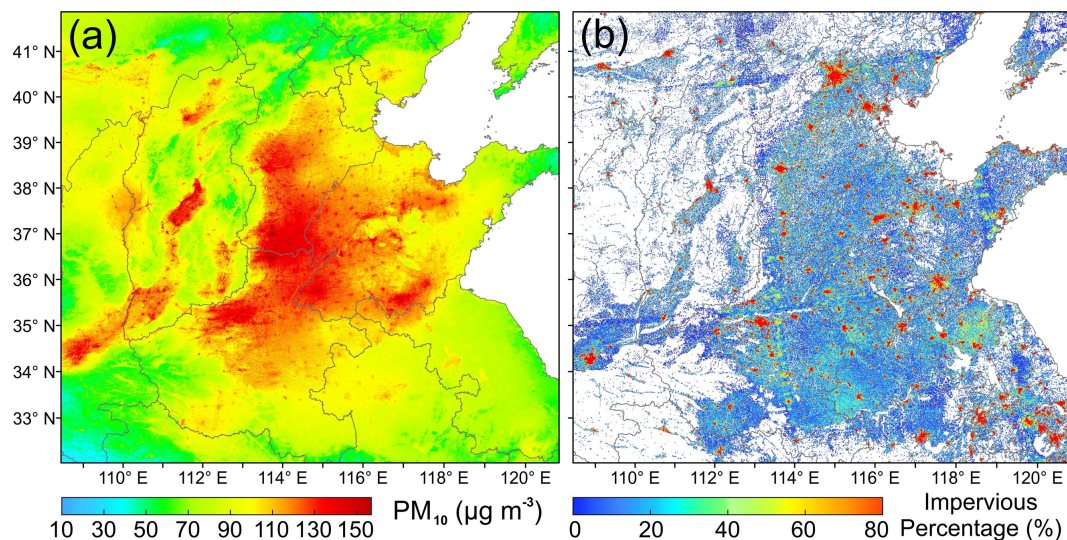
498 In Figure 8 we compared the spatial distribution of PM_{2.5} that was reconstructed by different
 499 datasets. Compared to LGHAP and TAP, PM_{2.5} data from CHAP are not gap free since the spatial
 500 coverage is determined by the AOD data coverage in the MAIAC product. Compared to TAP, LGHAP
 501 PM_{2.5} data have a finer resolution (1 km versus 10 km), enabling us to examine PM_{2.5} variations in
 502 space with more details. Overall, LGHAP has a better performance in reconstructing PM_{2.5} spatial
 503 distributions than the other two datasets. Reasons could be attributed to the following two aspects.
 504 Firstly, *in situ* PM_{2.5} measurements were fused with gridded PM_{2.5} estimations using the OI method
 505 when generating the final PM_{2.5} product in LGHAP. This can help correct modeling biases in original
 506 PM_{2.5} estimations. Secondly, a set of satellite-based AOD retrievals were incorporated when

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



510
 511 **Figure 8.** Comparison of PM_{2.5} distribution reconstructed by different PM_{2.5} concentration datasets.
 512 From the left to right, it shows in situ PM_{2.5} concentration measurements, CHAP, TAP, and LGHAP,
 513 respectively.

514
 515 To illustrate the fine resolution of LGHAP dataset, we compared the annual mean PM₁₀
 516 concentration in 2019 with the proportion of impervious surface that was derived from 30-m resolution
 517 land cover data in eastern China. As shown in Figure 9, the finer resolution of LGHAP dataset enables
 518 us to easily recognize the “hot spot” regions with high PM₁₀ loading. By referring to the impervious
 519 surface distribution on the right, we found that these hot spots are mainly over cities and towns,
 520 indicative of the presence of pollution island in urban regions. Owing to the involvement of such high-
 521 resolution datasets, the spatial details of PM_{2.5} and PM₁₀ can be then well recognized in LGHAP. The
 522 finer spatial resolution advantage of the LGHAP dataset can be also demonstrated by comparisons of
 523 spatial distribution of annual mean PM_{2.5} concentration that was revealed by four different datasets
 524 shown in Figure S6.



525

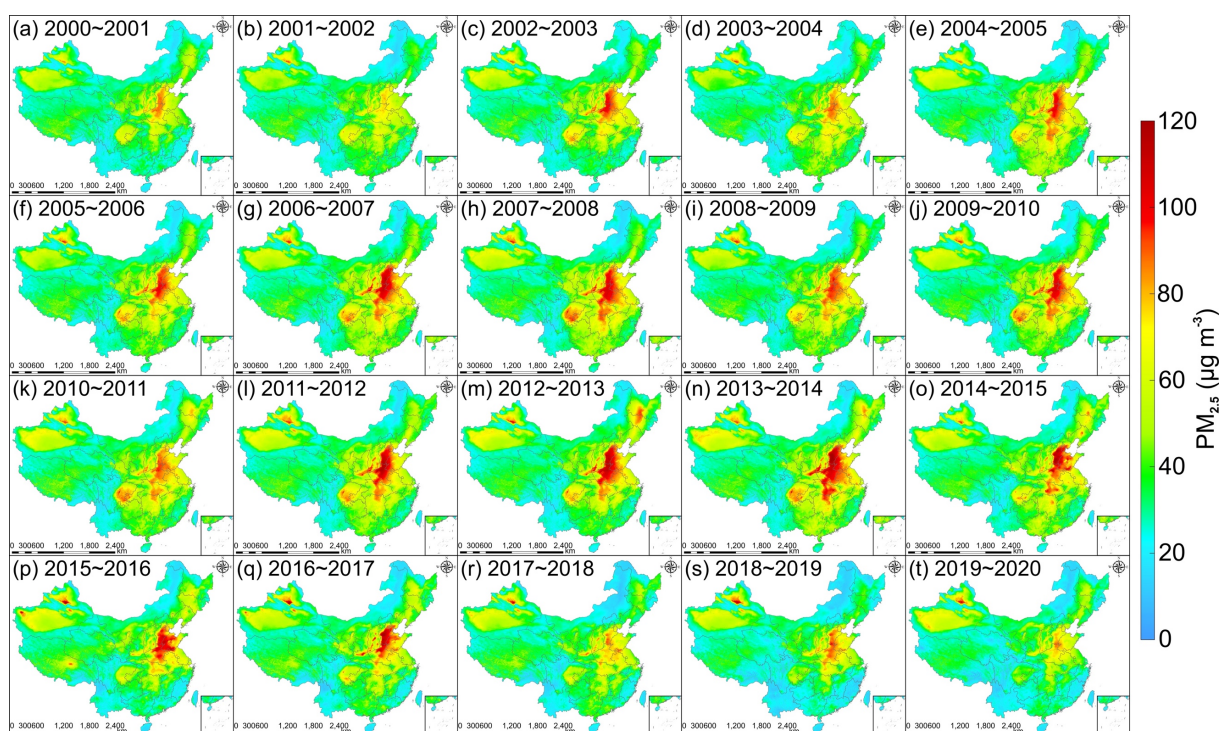
526 **Figure 9.** Comparison of annual mean PM₁₀ concentration with the proportion of areas covered by
 527 impervious surface in eastern China.

528 **4.3 Long-term trends of haze pollution in China from 2000 to 2020**

529 The aerosol pollution trends in China can be better examined by taking advantage of LGHAP
 530 dataset given long temporal coverage, gap free and high-resolution superiorities. Severe haze
 531 pollutions such as PM_{2.5} are oftentimes observed during winter half year (September–February). In
 532 this study, we first calculated mean PM_{2.5} concentration in China during winter half year from 2000 to
 533 2020. As shown in Figure 10, severe haze pollution events were mainly observed in North China during
 534 the wintertime, especially over the adjacent region in Hebei-Shandong-Henan provinces. In addition,
 535 Sichuan basin and Fenwei plain also suffered from severe haze pollution. Temporally, severe haze
 536 pollution events occurred mainly from the late 2002 to early 2017, which were significantly reduced
 537 after 2017. Similar pattern can be also inferred from PM₁₀ concentration distributions shown in Figure
 538 S7.

539 Figure 11 shows the temporal variations of the proportion of land areas covered by PM_{2.5}
 540 concentration exceeding 35 µg m⁻³ (the national ambient air quality standard for 24-hour PM_{2.5}
 541 concentration given in GB 3095-2012). As shown in Figure 11a, severe PM_{2.5} pollution occurred
 542 mainly during the wintertime in China, as more than one-third land areas (indicated by the blue lines)
 543 were exposed to unhealthy PM_{2.5} pollutants. Meanwhile, an apparent inflection was observed in 2007,
 544 after which the number of episode days decreased drastically at more than one-third land area covered
 545 by PM_{2.5} concentration exceeding 35 µg m⁻³. According to the proportion of land area covered with

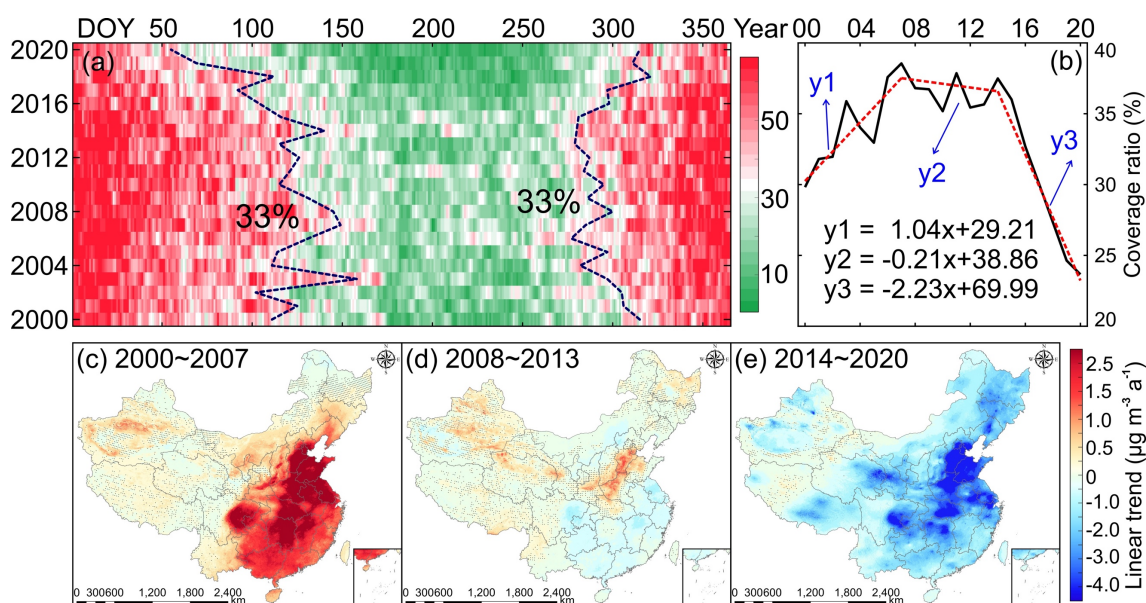
546 annual mean $\text{PM}_{2.5}$ concentration greater than $35 \mu\text{g m}^{-3}$, the variation of haze pollution in China can
 547 be generally divided into three different periods during the past two-decades (Figure 11b). As indicated,
 548 an increasing trend was observed from 2000 to 2007, during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



553
 554 **Figure 10.** Spatial distribution of mean $\text{PM}_{2.5}$ concentration from LGHAP during winter half year
 555 (September–February) from 2000 to 2020 in China.

556
 557 Figure 11c–e presents the linear trend of $\text{PM}_{2.5}$ concentration during these three specific periods,
 558 from which we observed that significant $\text{PM}_{2.5}$ variations occurred mainly over eastern part of the
 559 country where resides two-thirds of the population. A near ubiquitous $\text{PM}_{2.5}$ increasing trend was
 560 observed during 2000–2007, with significant increase ($>1.0 \mu\text{g m}^{-3} \text{ a}^{-1}$) mainly observed in eastern
 561 China. During the second period, $\text{PM}_{2.5}$ concentration over most regions shows a small decreasing
 562 trend except in the Ji-Lu-Yu region where an increasing trend was still observed. Apparent decreasing
 563 trend was observed over most parts of the country after 2014, indicative of significant reductions in
 564 $\text{PM}_{2.5}$ loading across China. This trend distribution is in line with our previous finding that was derived

565 using the annual mean PM_{2.5} concentration dataset generated by the Dalhousie University (Bai et al.,
 566 2019b). However, differences were still observed in terms of the regions where significant decreasing
 567 trends were present. Most significant decreasing trends were mainly observed in Sichuan basin and
 568 Pearl River Delta in the previous study. However, regions with drastic PM_{2.5} decrease were found
 569 mainly in the North China where severe haze pollution events were oftentimes reported. Similar
 570 variation patterns can be also inferred from PM₁₀ (Figure S8) and AOD (Figure S9). Overall, the
 571 LGHAP dataset provides us a gridded perspective to better examine long-term variations of haze
 572 pollution in China during the past two decades.

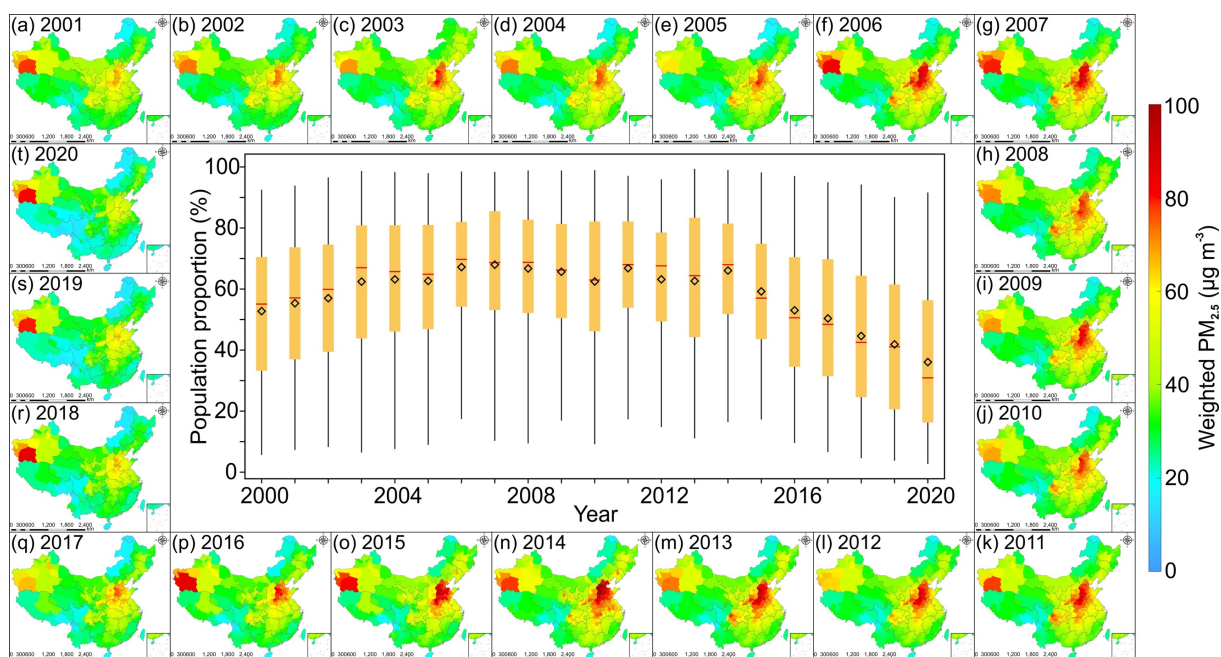


573
 574 **Figure 11.** Temporal variations of the proportion of land areas covered with PM_{2.5} concentration
 575 exceeding 35 µg m⁻³ and PM_{2.5} trends during three different periods. (a) Temporal variations of the
 576 land coverage ratio with daily PM_{2.5} concentration exceeding 35 µg m⁻³ from 2000 to 2000. (b) same
 577 as (a) but for annual mean PM_{2.5} concentration. (c–e) PM_{2.5} trends during periods of 2000–2007,
 578 2008–2013, and 2014–2020. The dotted regions imply trend estimations are statistically insignificant at the
 579 95% confidence interval.

580 4.4 Population exposure to PM_{2.5} pollution in China

581 By taking advantage of fine resolution LGHAP PM_{2.5} concentration and gridded population data,
 582 population exposure to PM_{2.5} pollution across China over the past two decades were estimated. Figure
 583 12 shows the spatial distribution of population weighted PM_{2.5} concentration and the proportion of
 584 population exposed to PM_{2.5} concentration greater than 35 µg m⁻³. As shown, spatial distribution of

585 population weighted $PM_{2.5}$ concentration resembles the spatial pattern of annual mean $PM_{2.5}$
 586 concentration, with high values observed mainly in eastern and central China as well as northwest
 587 China. Nonetheless, $PM_{2.5}$ sources in these two areas could be different. In northwest China, natural
 588 emissions could be the dominant source since very limited population resides there. In contrast, most
 589 population lives in eastern and central China with highly developed economy, and anthropogenic
 590 emissions thus might play more important roles in $PM_{2.5}$ formation (Xin et al., 2015; Yang et al., 2011).
 591 In regard to the proportion of population exposed to the ambient with $PM_{2.5}$ concentration greater than
 592 $35 \mu g m^{-3}$, we observed that the annual mean population ratio exposure to unhealthy $PM_{2.5}$ increased
 593 gradually from 50.60% in 2000 to 65.72% in 2007. During 2007–2014, the ratio varied with small
 594 changes (<5%), whereas a drastic decline was observed after 2014, with the annual mean proportion
 595 of population exposed to unhealthy $PM_{2.5}$ was reduced from 63.81% in 2014 to 34.03% in 2020, even
 596 though the total population was increased from 1.37 billion to 1.41 billion during the synchronous
 597 period. Nonetheless, more than one-third population was still exposed to unhealthy $PM_{2.5}$, highlighting
 598 the requirement of further emission reduction actions to manage haze pollutions in China.



599 **Figure 12.** Spatial distribution of population weighted $PM_{2.5}$ concentration and the proportion of
 600 population exposed to $PM_{2.5}$ concentration greater than $35 \mu g m^{-3}$. Annual and daily LGHAP $PM_{2.5}$
 601 concentration data were used for the calculation of weighted $PM_{2.5}$ and the proportion of population
 602 exposure, respectively. The diamond and red line indicate the annual mean and median population
 603 proportion, respectively.
 604

605 **5 Data availability**

606 The LGHAP dataset, consisting of gap free AOD, PM_{2.5}, and PM₁₀ concentration with daily 1-
607 km resolution from 2000 to 2020, are all publicly accessible. All data were provided in the NetCDF
608 format and data in each individual year were archived in a zip file. For AOD, the dataset has a disk
609 storage size of near 27 GB in total, which is available at <https://doi.org/10.5281/zenodo.5652257> (Bai et
610 al., 2021a). PM_{2.5} (38 GB) and PM₁₀ (48 GB) concentration data can be acquired from
611 <https://doi.org/10.5281/zenodo.5652265> (Bai et al., 2021b) and <https://doi.org/10.5281/zenodo.5652263> (Bai
612 et al., 2021c), respectively. Additionally, monthly and annual mean datasets were also provided, which
613 is publicly available at <https://doi.org/10.5281/zenodo.5655797> (Bai et al., 2021d) and
614 <https://doi.org/10.5281/zenodo.5655807> (Bai et al., 2021e), respectively. In addition to these datasets,
615 Python, Matlab, R, and IDL codes that can be used to read and visualize these data were provided as
616 well.

617 **6 Conclusion**

618 In this study, a big data analytics method was developed for generating a LGHAP dataset to
619 advance research in earth system science and environment management. With integrative efforts of
620 fusing AOD features extracted from a set of AOD data tensors and knowledge transfer in statistical
621 data mining from diverse air quality indicators, a LGHAP aerosol dataset providing 21-year-long
622 (2000–2020) gap-free AOD, PM_{2.5}, and PM₁₀ concentration data with daily 1-km resolution in China,
623 was generated. Gap-filled AOD imageries were firstly generated by reconstructing AOD distribution
624 in AOD_{Terra} via synergistically fusing AOD features recognized from diversified satellites and
625 numerical models as well as *in situ* data through tensor completion. Compared to ground-based AOD
626 measurements, the gap-filled AOD data exhibit a satisfying prediction accuracy and good performance
627 in delineating AOD variations over space and time. To our knowledge, this is the first thrust of
628 generating long-term high-resolution AOD dataset with gap free nature in China.

629 PM_{2.5} and PM₁₀ concentration data were then estimated using an ensemble learning approach by
630 taking advantage of the generated gap-free AOD imageries. Ground validation results also indicate
631 good accuracies of these two gridded products, showing a comparable bias level with many previous
632 studies. Compared with other open access daily PM_{2.5} concentration datasets, the LGHAP PM_{2.5}
633 dataset performs well due to the vantage of having gap free and fine resolution products. With this gap

634 free and high-resolution dataset, the long-term variation trend of haze pollution in China over the past
635 two decades was examined, and apparent inflections were observed in 2007 and 2014, at which PM_{2.5}
636 concentration was found to turn from an increasing path to decreasing in 2007 with a more drastic
637 decline observed starting from 2014. Moreover, the LGHAP dataset provides us a gridded perspective
638 to assess two-decade long population exposure to PM_{2.5} pollution in China. In spite of a drastic decline
639 in population exposure, there are still more than one-third population exposed to unhealthy PM_{2.5}
640 pollutants, highlighting the requirement of long-lasting actions to continue PM_{2.5} related emission
641 reduction.

642 Overall, these three gridded LGHAP aerosol products provide a long-term perspective on aerosol
643 changes over different regions of China, and users are encouraged to use the LGHAP dataset to assess
644 aerosol impacts on public health, air quality, climate, and ecosystem. The dataset has been publicly
645 released online and is freely accessible via the links provided in Section 5. Global scale dataset is on
646 the track and will be released to the public soon.

647 **Author contributions**

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

653 **Competing interests**

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

655 **Acknowledgments**

656 The authors are grateful to the editor and two anonymous referees for their constructive comments and
657 suggestions in improving this manuscript. Also, the authors would like to thank all organizations and
658 groups for providing essential datasets that were used in this study. The MAIAC AOD was acquired
659 from <https://lpdaac.usgs.gov/products/mcd19a2v006/>. The MISR AOD was acquired from
660 <https://asdc.larc.nasa.gov/project/MISR>. The VIIRS AOD was acquired from

661 <https://earthdata.nasa.gov/earth-observation-data/near-real-time/download-nrt-data/viirs-nrt>. The
662 AATSR AOD was acquired from <https://climate.esa.int/en/projects/aerosol/data/>. The POLDER AOD
663 was acquired from <https://www.grasp-open.com/products/polder-data-release/>. The aerosol
664 diagnostics including AOD and aerosol components from MERRA-2 were acquired from
665 https://disc.gsfc.nasa.gov/datasets/M2T1NXAER_5.12.4/summary?keywords=MERRA2. AOD from
666 AERONET was acquired from https://aeronet.gsfc.nasa.gov/new_web/aerosols.html. Meteorological
667 factors were retrieved from the latest ERA-5 reanalysis and can be reached at
668 <https://cds.climate.copernicus.eu/cdsapp#!/dataset/reanalysis-era5-single-levels?tab=overview>.
669 Atmospheric visibility data were acquired from the national meteorological information center at
670 <http://data.cma.cn/en>. Ground-based air pollutants concentration was acquired from
671 <https://air.cnemc.cn:18007/>. Gridded Population data were acquired from <https://www.worldpop.org/>
672 while DEM was acquired from <https://www.resdc.cn/>. Monthly NDVI data were acquired from
673 <https://lpdaac.usgs.gov/products/mod13a3v061/>. Land cover data were acquired from
674 [http://www.globallandcover.com/defaults.html?src=/Scripts/map/defaults/browse.html&head=browse](http://www.globallandcover.com/defaults.html?src=/Scripts/map/defaults/browse.html&head=browse&type=data)
675 [e&type=data](https://zenodo.org/record/4417810#.YSxD844zYuW) and <https://zenodo.org/record/4417810#.YSxD844zYuW>.

676 **Financial support**

677 This study was supported by the National Natural Science Foundation of China (grants 42171309 and
678 41701413), and the Shanghai Committee of Science and Technology (grant 20ZR1415900).

679 **References**

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

683 Bai, K., Li, K., Chang, N.-B. and Gao, W.: Advancing the prediction accuracy of satellite-based
684 PM_{2.5} concentration mapping: A perspective of data mining through in situ PM_{2.5} measurements,
685 *Environ. Pollut.*, 254, 113047, doi:10.1016/j.envpol.2019.113047, 2019a.

686 Bai, K., Ma, M., Chang, N.-B. and Gao, W.: Spatiotemporal trend analysis for fine particulate
687 matter concentrations in China using high-resolution satellite-derived and ground-measured PM_{2.5}
688 data, *J. Environ. Manage.*, 233, 530–542, doi:10.1016/j.jenvman.2018.12.071, 2019b.

689 Bai, K., Li, K., Wu, C., Chang, N.-B. and Guo, J.: A homogenized daily in situ PM_{2.5}
690 concentration dataset from the national air quality monitoring network in China, *Earth Syst. Sci. Data*,
691 12(4), 3067–3080, doi:10.5194/essd-12-3067-2020, 2020a.

692 Bai, K., Li, K., Guo, J., Yang, Y. and Chang, N.-B.: Filling the gaps of in situ hourly PM_{2.5}
693 concentration data with the aid of empirical orthogonal function analysis constrained by diurnal cycles,
694 *Atmos. Meas. Tech.*, 13(3), 1213–1226, doi:10.5194/amt-13-1213-2020, 2020b.

695 Bai, K., Li, K., Guo, J. and Chang, N.-B.: Multiscale and multisource data fusion for full-coverage
696 PM_{2.5} concentration mapping: Can spatial pattern recognition come with modeling accuracy?, *ISPRS*
697 *J. Photogramm. Remote Sens.*, 184, 31–44, doi: 10.1016/j.isprsjprs.2021.12.002, 2022.

698 Bai, K., Li, K. Tan, Z., Han, D., and Guo, J.: Daily 1-km gap-free AOD grids in China, v1 (2000–
699 2020) [data set], <https://doi.org/10.5281/zenodo.5652257>, 2021a.

700 Bai, K., Li, K. Tan, Z., Han, D., and Guo, J.: Daily 1-km gap-free PM_{2.5} grids in China, v1 (2000–
701 2020) [data set], <https://doi.org/10.5281/zenodo.5652265>, 2021b.

702 Bai, K., Li, K. Tan, Z., Han, D., and Guo, J.: Daily 1-km gap-free PM₁₀ grids in China, v1 (2000–
703 2020) [data set], <https://doi.org/10.5281/zenodo.5652263>, 2021c.

704 Bai, K., Li, K. Tan, Z., Han, D., and Guo, J.: Monthly averaged 1-km gap-free AOD, PM_{2.5} and
705 PM₁₀ grids in China, v1 (2000–2020) [data set], <https://doi.org/10.5281/zenodo.5655797>, 2021d.

706 Bai, K., Li, K. Tan, Z., Han, D., and Guo, J.: Annual mean 1-km gap-free AOD, PM_{2.5} and PM₁₀
707 grids in China, v1 (2000–2020) [data set], <https://doi.org/10.5281/zenodo.5655807>, 2021e.

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

709 Datasets, *J. Atmos. Ocean. Technol.*, 20(12), 1839–1856, doi:10.1175/1520-
710 0426(2003)020<1839:ECADFF>2.0.CO;2, 2003.

711 Bi, J., Belle, J. H., Wang, Y., Lyapustin, A. I., Wildani, A. and Liu, Y.: Impacts of snow and cloud
712 covers on satellite-derived PM_{2.5} levels, *Remote Sens. Environ.*, 221(October), 665–674,
713 doi:10.1016/j.rse.2018.12.002, 2018.

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

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

720 Chen, G., Li, S., Knibbs, L. D., Hamm, N. A. S., Cao, W., Li, T., Guo, J., Ren, H., Abramson, M.
721 J. and Guo, Y.: A machine learning method to estimate PM_{2.5} concentrations across China with
722 remote sensing, meteorological and land use information, *Sci. Total Environ.*, 636, 52–60,
723 doi:10.1016/j.scitotenv.2018.04.251, 2018.

724 Chen, J., Ban, Y., and Li, S.: China: Open access to Earth land-cover map, *Nature*, 514(7523):
725 434-434, doi:10.1038/514434c, 2014.

726 Di, Q., Amini, H., Shi, L., Kloog, I., Silvern, R., Kelly, J., Sabath, M. B., Choirat, C., Koutrakis,
727 P., Lyapustin, A., Wang, Y., Mickley, L. J. and Schwartz, J.: An ensemble-based model of PM_{2.5}
728 concentration across the contiguous United States with high spatiotemporal resolution, *Environ. Int.*,
729 130, 104909, doi:10.1016/j.envint.2019.104909, 2019.

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

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

738 Fang, X., Zou, B., Liu, X., Sternberg, T. and Zhai, L.: Satellite-based ground PM_{2.5} estimation

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

741 Fuzzi, S., Baltensperger, U., Carslaw, K., Decesari, S., Denier van der Gon, H., Facchini, M. C.,
742 Fowler, D., Koren, I., Langford, B., Lohmann, U., Nemitz, E., Pandis, S., Riipinen, I., Rudich, Y.,
743 Schaap, M., Slowik, J. G., Spracklen, D. V., Vignati, E., Wild, M., Williams, M. and Gilardoni, S.:
744 Particulate matter, air quality and climate: lessons learned and future needs, *Atmos. Chem. Phys.*,
745 15(14), 8217–8299, doi:10.5194/acp-15-8217-2015, 2015.

746 Gao, M., Beig, G., Song, S., Zhang, H., Hu, J., Ying, Q., Liang, F., Liu, Y., Wang, H., Lu, X.,
747 Zhu, T., Carmichael, G. R., Nielsen, C. P. and McElroy, M. B.: The impact of power generation
748 emissions on ambient PM_{2.5} pollution and human health in China and India, *Environ. Int.*,
749 121(August), 250–259, doi:10.1016/j.envint.2018.09.015, 2018.

750 Geng, G., Xiao, Q., Liu, S., Liu, X., Cheng, J., Zheng, Y., Xue, T., Tong, D., Zheng, B., Peng, Y.,
751 Huang, X., He, K. and Zhang, Q.: Tracking Air Pollution in China: Near Real-Time PM_{2.5} Retrievals
752 from Multisource Data Fusion, *Environ. Sci. Technol.*, 55(17), 12106–12115,
753 doi:10.1021/acs.est.1c01863, 2021.

754 Goldberg, D. L., Gupta, P., Wang, K., Jena, C., Zhang, Y., Lu, Z. and Streets, D. G.: Using gap-
755 filled MAIAC AOD and WRF-Chem to estimate daily PM_{2.5} concentrations at 1 km resolution in the
756 Eastern United States, *Atmos. Environ.*, 199(November 2018), 443–452,
757 doi:10.1016/j.atmosenv.2018.11.049, 2019.

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

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

765 Huang, C., Hu, J., Xue, T., Xu, H. and Wang, M.: High-Resolution Spatiotemporal Modeling for
766 Ambient PM_{2.5} Exposure Assessment in China from 2013 to 2019, *Environ. Sci. Technol.*, 55(3),
767 2152–2162, doi:10.1021/acs.est.0c05815, 2021.

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

769 500, doi:10.1137/07070111X, 2009.

770 de Leeuw, G., Sogacheva, L., Rodriguez, E., Kourtidis, K., Georgoulas, A. K., Alexandri, G.,
771 Amiridis, V., Proestakis, E., Marinou, E., Xue, Y. and van der A, R.: Two decades of satellite
772 observations of AOD over mainland China using ATSR-2, AATSR and MODIS/Terra: data set
773 evaluation and large-scale patterns, *Atmos. Chem. Phys.*, 18(3), 1573–1592, doi:10.5194/acp-18-
774 1573-2018, 2018.

775 Li, J., Li, C. and Zhao, C.: Different trends in extreme and median surface aerosol extinction
776 coefficients over China inferred from quality-controlled visibility data, *Atmos. Chem. Phys.*, 18(5),
777 3289–3298, doi:10.5194/acp-18-3289-2018, 2018a.

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

781 Li, K., Bai, K., Li, Z., Guo, J. and Chang, N.-B.: Synergistic Data Fusion of Multimodal AOD and
782 Air Quality Data for Near Real-Time Full Coverage Air Pollution Assessment, *J. Environ. Manage.*,
783 302, 114121, doi: 10.1016/j.jenvman.2021.114121, 2022.

784 Li, Z., Zhang, Y., Shao, J., Li, B., Hong, J., Liu, D., Li, D., Wei, P., Li, W., Li, L., Zhang, F., Guo,
785 J., Deng, Q., Wang, B., Cui, C., Zhang, W., Wang, Z., Lv, Y., Xu, H., Chen, X., Li, L. and Qie, L.:
786 Remote sensing of atmospheric particulate mass of dry PM_{2.5} near the ground: Method validation
787 using ground-based measurements, *Remote Sens. Environ.*, 173, 59–68, doi:10.1016/j.rse.2015.11.019,
788 2016.

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

794 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,
795 X. and Yu, Q.: Estimation of long-term population exposure to PM_{2.5} for dense urban areas using 1-
796 km MODIS data, *Remote Sens. Environ.*, 179, 13–22, doi:10.1016/j.rse.2016.03.023, 2016.

797 Liu, M., Bi, J. and Ma, Z.: Visibility-Based PM_{2.5} Concentrations in China: 1957–1964 and
798 1973–2014, *Environ. Sci. Technol.*, 51(22), 13161–13169, doi:10.1021/acs.est.7b03468, 2017.

799 Liu, Y., Paciorek, C. J. and Koutrakis, P.: Estimating Regional Spatial and Temporal Variability
800 of PM_{2.5} Concentrations Using Satellite Data, Meteorology, and Land Use Information, *Environ.*
801 *Health Perspect.*, 117(6), 886–892, doi:10.1289/ehp.0800123, 2009.

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

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

807 Lyu, B., Hu, Y., Zhang, W., Du, Y., Luo, B., Sun, X., Sun, Z., Deng, Z., Wang, X., Liu, J., Wang,
808 X. and Russell, A. G.: Fusion Method Combining Ground-Level Observations with Chemical
809 Transport Model Predictions Using an Ensemble Deep Learning Framework: Application in China to
810 Estimate Spatiotemporally-Resolved PM_{2.5} Exposure Fields in 2014-2017, *Environ. Sci. Technol.*,
811 53(13), 7306–7315, doi:10.1021/acs.est.9b01117, 2019.

812 Ma, Z., Hu, X., Sayer, A. M., Levy, R., Zhang, Q., Xue, Y., Tong, S., Bi, J., Huang, L. and Liu,
813 Y.: Satellite-based spatiotemporal trends in PM_{2.5} concentrations: China, 2004-2013, *Environ. Health*
814 *Perspect.*, 124(2), 184–192, doi:10.1289/ehp.1409481, 2016.

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

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

823 Sidiropoulos, N. D., De Lathauwer, L., Fu, X., Huang, K., Papalexakis, E. E. and Faloutsos, C.:
824 Tensor Decomposition for Signal Processing and Machine Learning, *IEEE Trans. Signal Process.*,
825 65(13), 3551–3582, doi:10.1109/TSP.2017.2690524, 2017.

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

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

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

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

838 Tang, Q., Bo, Y. and Zhu, Y.: Spatiotemporal fusion of multiple-satellite aerosol optical depth
839 (AOD) products using Bayesian maximum entropy method, *J. Geophys. Res. Atmos.*, 121(8), 4034–
840 4048, doi:10.1002/2015JD024571, 2016.

841 Tucker, L. R.: Some mathematical notes on three-mode factor analysis, *Psychometrika*, 31(3),
842 279–311, doi:10.1007/BF02289464, 1966.

843 Wang, B., Yuan, Q., Yang, Q., Zhu, L., Li, T. and Zhang, L.: Estimate hourly PM_{2.5}
844 concentrations from Himawari-8 TOA reflectance directly using geo-intelligent long short-term
845 memory network, *Environ. Pollut.*, 271, 116327, doi:10.1016/j.envpol.2020.116327, 2021a.

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

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

851 Wang, Y., Yuan, Q., Li, T., Tan, S. and Zhang, L.: Full-coverage spatiotemporal mapping of
852 ambient PM_{2.5} and PM₁₀ over China from Sentinel-5P and assimilated datasets: Considering the
853 precursors and chemical compositions, *Sci. Total Environ.*, 793, 148535,
854 doi:10.1016/j.scitotenv.2021.148535, 2021b.

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

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

859 km-resolution high-quality PM_{2.5} data records from 2000 to 2018 in China: spatiotemporal variations
860 and policy implications, *Remote Sens. Environ.*, 252(January 2020), 112136,
861 doi:10.1016/j.rse.2020.112136, 2021a.

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

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

868 Xiao, Q., Zhang, H., Choi, M., Li, S., Kondragunta, S., Kim, J., Holben, B., Levy, R. C. and Liu,
869 Y.: Evaluation of VIIRS, GOCI, and MODIS Collection 6 AOD retrievals against ground
870 sunphotometer observations over East Asia, *Atmos. Chem. Phys.*, 16(3), 1255–1269, doi:10.5194/acp-
871 16-1255-2016, 2016.

872 Xiao, Q., Wang, Y., Chang, H. H., Meng, X., Geng, G., Lyapustin, A. and Liu, Y.: Full-coverage
873 high-resolution daily PM_{2.5} estimation using MAIAC AOD in the Yangtze River Delta of China,
874 *Remote Sens. Environ.*, 199(May), 437–446, doi:10.1016/j.rse.2017.07.023, 2017.

875 Xiao, Q., Chang, H. H., Geng, G. and Liu, Y.: An Ensemble Machine-Learning Model to Predict
876 Historical PM_{2.5} Concentrations in China from Satellite Data, *Environ. Sci. Technol.*,
877 doi:10.1021/acs.est.8b02917, 2018.

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

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

885 Xue, T., Zheng, Y., Tong, D., Zheng, B., Li, X., Zhu, T. and Zhang, Q.: Spatiotemporal continuous
886 estimates of PM_{2.5} concentrations in China, 2000–2016: A machine learning method with inputs from
887 satellites, chemical transport model, and ground observations, *Environ. Int.*, 123(December 2018),
888 345–357, doi:10.1016/j.envint.2018.11.075, 2019.

889 Yang, F., Tan, J., Zhao, Q., Du, Z., He, K., Ma, Y., Duan, F., Chen, G. and Zhao, Q.:
890 Characteristics of PM_{2.5} speciation in representative megacities and across China, *Atmos. Chem.*
891 *Phys.*, 11(11), 5207–5219, doi:10.5194/acp-11-5207-2011, 2011.

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

894 Yang Y., Zheng Z., Yim S.H.L., Roth M., Ren G., Gao Z., Wang T., Li Q., Shi C., Ning G. and
895 Li Y.B.: PM_{2.5} Pollution Modulates Wintertime Urban-Heat-Island Intensity in the Beijing-Tianjin-
896 Hebei Megalopolis, China. *Geophys. Res. Lett.*, 47(1), e2019GL084288, doi:10.1029/2019gl084288,
897 2020.

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

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

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

907 Zheng, Z., Ren, G., Wang, H., Dou, J., Gao, Z., Duan, C., Li, Y. and Ngarukiyimana, J., Zhao,
908 C., Cao, C., Jiang, M., and Yang, Y.: Relationship between Fine Particle Pollution and the Urban Heat
909 Island in Beijing, China: Observational Evidence, *Boundary-Layer Meteorol.*, 169(1), 93-113, doi:
910 10.1007/s10546-018-0362-6, 2018.

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

914

915