

# Global aerosol typing classification using a new hybrid algorithm utilizing Aerosol Robotic Network data

Xiaoli Wei<sup>1,2</sup>, Qian Cui<sup>5</sup>, Leiming Ma<sup>1</sup>, Feng Zhang<sup>2,3</sup>, Wenwen Li<sup>2,3</sup>, Peng Liu<sup>4</sup>

<sup>1</sup> Shanghai Meteorological Service 200030, China;

<sup>2</sup> Shanghai Qi Zhi Institute, Shanghai, 200232, China;

<sup>3</sup> Department of Atmospheric and Oceanic Sciences & Institute of Atmospheric Sciences, Fudan University, Shanghai, 200438, China;

<sup>4</sup> School of Atmospheric Science, Nanjing University of Information Science and Technology, Nanjing 210044, China;

<sup>5</sup> Caidian Meteorological Service, Wuhan, 430000, China

Correspondence to: Feng Zhang (fengzhang@fudan.edu.cn)

## Abstract

Aerosols have great uncertainty owing to the complex changes in their composition in different regions. The radiation properties of different aerosol types differ considerably and are vital in studying aerosol regional and/or global climate effects. Traditional aerosol-type identification algorithms, generally based on cluster or empirical analysis methods, are often inaccurate and time-consuming. In response, our study aimed to develop a new aerosol-type classification model using an innovative hybrid algorithm to improve the precision and efficiency of aerosol-type identification. This novel algorithm incorporates an optical database, constructed using the Mie scattering model, and employs a random forest algorithm to classify different aerosol types based on the optical data from the database. The complex refractive index was used as a baseline to assess the performance of our hybrid algorithm against the traditional Gaussian kernel density clustering method for aerosol type identification. The hybrid algorithm demonstrated impressive consistency rates of 90%, 85%, 84%, 84%, and 100% for dust, mixed-coarse, mixed-fine, urban/industrial, and biomass burning aerosols, respectively. Moreover, it achieved remarkable precision, with F-score and accuracy scores of 95%, 89%, 91%, and 89%. Lastly, a global map of aerosol types was generated using the new hybrid algorithm to characterize aerosol types across the five continents. This study utilizing a novel

31 approach for the classification of aerosol will help improve the accuracy of aerosol  
32 inversion and determine the sources of aerosol pollution.

33 **Keywords:** Aerosol typing classification, Hybrid algorithm, Complex refractive index,  
34 AERONET

## 35 **1. Introduction**

36 Atmospheric aerosols are tiny solid or liquid particles suspended in the  
37 atmosphere. Aerosols indirectly affect the energy budget and water cycle of the earth's  
38 gas system by absorbing and scattering solar radiation or by changing the optical  
39 properties and life cycle of the cloud as condensation nuclei of cloud droplets  
40 (Redemann et al. 2000; Ramanathan et al. 2001). Additionally, desert dust, biomass  
41 smog, and anthropogenic emissions of air pollutants can affect visibility, air quality,  
42 and human health (Tong et al., 2017; Siomos et al., 2020). Evaluating the impact of  
43 aerosols on radiative transfer is complex, primarily because of the uncertainty of  
44 radiative forcing caused by the high spatiotemporal dynamic variation of aerosol  
45 optical and physical characteristics in different regions (Kaskaoutis et al., 2011; Che et  
46 al., 2018; Elham et al., 2023). The aerosol type embodies the long-term average  
47 physicochemical properties of aerosols in a certain area (Kiehl & Briegleb, 1993; Lu  
48 et al., 2023). Therefore, accurate identification of aerosol types can drive the study of  
49 the climatic effects of aerosols, tracking and control of environmental pollution  
50 sources, and precision of radiation transmission models.

51 Aerosol types are defined based on the radiation properties of different aerosol  
52 types owing to the large variation in their optical, physical, and chemical properties.  
53 Currently, aerosol types are classified by two ways by using the traditional clustering  
54 algorithms (Kumar et al., 2018). First, based on different sources and properties at  
55 different observation points worldwide, aerosols are classified as follows: dust  
56 aerosols from deserts, biomass combustion aerosols from forests or grasslands, and  
57 urban/industrial (U/I) aerosols from fuel combustion in densely populated urban areas  
58 (Dubovik et al., 2002; Pawar et al., 2015; Yousefi et al., 2020). Second, based on the

59 size of the radiation absorption rate, aerosols into four categories: carbonaceous (fine-  
60 absorbing mode), soil dust (coarse absorption mode), sulfates (nonabsorbing fine-  
61 grained mode), and sea salt aerosols (nonabsorbing coarse-grained mode) (Levy et al.,  
62 2007). The first classification, widely used for aerosol retrieval and common in  
63 research, categorizes aerosol types based on optical properties observed at ground  
64 stations. This forms a two-dimensional identification space for clustering, while the  
65 second approach specifically subcategorizes anthropogenic aerosols. Many  
66 combinations of optical properties and parameters are available, such as  $EAE_{440-870nm}$   
67 (extinction angstrom exponent) vs.  $SSA_{440nm}$  (single-scattering albedo),  $AAE_{440-870nm}$   
68 (absorption angstrom exponent) vs.  $EAE_{440-870nm}$ ,  $AAE_{440-870nm}$  vs.  $FMF_{550nm}$  (fine  
69 mode fraction), and  $SSA_{440nm}$  vs.  $EAE_{440-870nm}$  (Lee et al., 2010; Shin et al., 2019; Choi,  
70 et al., 2021). Various studies have highlighted the importance of selecting appropriate  
71 aerosol properties for accurate aerosol type identification (Giles et al., 2012; Che et al.,  
72 2018).

73 Among the aerosol-type classification methodologies developed, those using  
74 threshold and empirical analyses have the greatest potential for large-area and fixed-  
75 period applications (Eck et al., 1999; Omar et al., 2005; Yang et al., 2009).  
76 Traditionally, the aerosol-type classification algorithm mainly distinguishes different  
77 aerosol types based on their optical properties and determines the threshold of their  
78 optical properties based on clustering. However, the composition of aerosols changes  
79 rapidly with time and location, owing to the combined influence of natural conditions  
80 and human activities (for example, tornadoes and various anthropogenic activities)  
81 (Sheridan et al., 2001). Unfortunately, determining aerosol types accurately and  
82 rapidly is a challenge when using traditional methods (Bahadur et al., 2012; Shin et al.,  
83 2019; Lin et al., 2021). Nevertheless, with advancements in data science, artificial  
84 intelligence techniques have aided the accurate and rapid recognition of different  
85 aerosol types.

86 Artificial intelligence algorithms can receive multiple aerosol characteristic  
87 parameters as input, thus preventing the sole reliance of aerosol classification on a  
88 limited number of features (Li et al., 2022; Wang et al., 2023). For example, Boselli

89 (2012) performed a k-means clustering analysis of single scattering albedo (SSA),  
90 aerosol optical depth (AOD), electrical asymmetry effect (EAE), and asymmetry  
91 parameter ( $g$ ) datasets for the central Mediterranean Sea for the classification of  
92 aerosol into four: dusty, continental, oceanic, or mixed aerosols. Nicolae (2018)  
93 developed a neural network algorithm to estimate the aerosol typing of Lidar data and  
94 Hamill (2016) introduced the Mahalanobis Distance for aerosol classification to  
95 determine a specific aerosol type for each reference cluster. Li (2022) generated  
96 spatial contiguous aerosol type map in China with an empirical aerosol type retrieval  
97 algorithm. Overall, limited information on the optical properties of aerosols can  
98 reasonably determine the type of aerosol (Hamill et al., 2016). However, some  
99 challenges remain in identifying aerosol types through machine learning. First, the  
100 amount of valid ground aerosol property data that can be used for training is less due  
101 to cloud removal and quality control. Second, the accuracy of machine learning  
102 depends on the labeled aerosol typing dataset, and finding a suitable classification  
103 method to classify the dataset is challenging. Third, evaluating the accuracy of the  
104 final trained model is also tedious (Zhang & Li, 2019; Siomos et al., 2020; Choi, et al.,  
105 2021a,b)

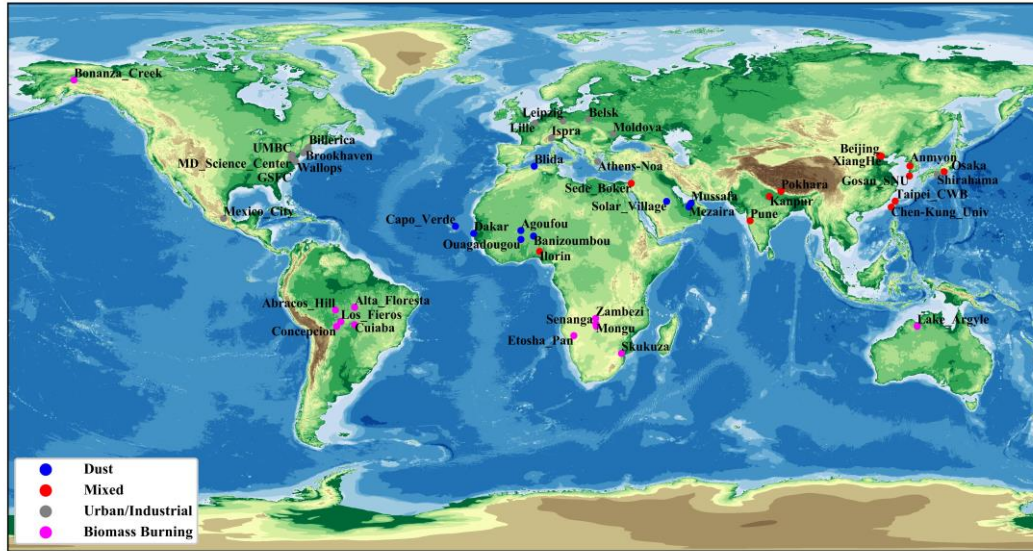
106 The traditional aerosol type identification methods are easily limited by time and  
107 space, and most of them only classify aerosol types using two optical property  
108 parameters, limiting the complete characterization of aerosols. Considering these  
109 limitations, we aimed to (1) develop a new algorithm that can accurately and quickly  
110 identify aerosol types to overcome existing problems such as low accuracy,  
111 insufficient data, and difficulty in setting labels; (2) investigate the characteristics of  
112 the regional spatial distribution of global aerosol types obtained using the new  
113 machine learning algorithms, considering the large regional differences in aerosol  
114 types. To achieve this, we propose a new aerosol-type classification algorithm based  
115 on a Gaussian cluster and random forest algorithm to generate an aerosol-typing map  
116 over several representative regions of the world.

## 117 2. Study area and data

118 Figure 1 illustrates the research area and the distribution of the Aerosol Robotic  
119 Network (AERONET) sites, strategically encompassing major global regions to  
120 validate the universality of the research algorithm. The study utilized 47 marked  
121 aerosol sites across five continents, leveraging them to train and validate the machine  
122 learning approach based on a comprehensive literature review. The 47 sites represent  
123 different aerosol-type properties of different aerosol source regions, including dust,  
124 mixed (mixed coarse and mixed fine aerosols), U/I, and biomass burning (BB)  
125 aerosols (Table 1 and Figure 1). Marine aerosols were not considered because their  
126 low optical thickness values (generally  $<0.4$ ) can result in a less valid data scale that  
127 would not meet the study requirements. Here, the aerosol source region refers to the  
128 area affected by one dominant emission source, where the aerosol types are fixed and  
129 not easily confused (Giles et al., 2012; Hamill et al., 2016). Table 2 presents the  
130 optical properties and microphysical characteristic parameters of aerosols at four  
131 bands of AERONET (440, 675, 870, and 1020 nm). These parameters were used to  
132 construct a database of SSA, AOD, and asymmetry parameters. Further, typical sites  
133 dominated by different aerosol types worldwide were selected for compositional  
134 analysis using the new model. The selected sites are distributed across different  
135 regions of the world and represent a specific aerosol-dominated type and aerosol  
136 source region.

137 For dust aerosols, five AERONET sites, namely Banizoumbou, Capo\_Verde,  
138 Dakar, and Ouagadougou in Africa and Solar\_Village in West Asia, influenced by the  
139 Saharan Desert, were considered. The Dakar and Capo\_Verde sites are located at the  
140 tip of the Capo\_Verde Peninsula—the westernmost part of Africa, bordering the  
141 Atlantic Ocean. Despite being oceanic, these two sites are dominated by dust aerosols  
142 influenced by aerosol plumes in the Saharan Desert. Meanwhile, the Banizoumbou  
143 and Ouagadougou are centrally located in Africa. Here, northeasterly winds in winter  
144 and northwesterly winds in summer transport Saharan Desert dust aerosols. For mixed  
145 aerosols, the AERONET sites Ilorin, Kanpur, Sede\_Boker, and XiangHe were

146 selected. For U/I aerosols, the AERONET sites GSFC, Ispra, Mexico\_City, and  
 147 Moldova were selected. Four AERONET sites, namely, Alta\_Floresta, Abracos\_Hill,  
 148 Lake\_Argyle, and Mongu, were selected as BB aerosol-dominant sites.



149  
 150 **Figure 1.** Study area and 47 AERONET sites selected by literature review.

151 **Table 1.** 47 AERONET sites selected by literature review.

Aerosol Type	Sites for Training	Sites for Testing
Dust	Agoufou, Capu_Verde, Dakar, Mezzaira, Mussafa, Ouagadougou	Banizoumbou, Solar_Village, Blida
Mixed	Anmyon, Beijing, Chen-Kung_Univ, Ilorin, Kanpur, Sede_Boker, Gosan_SUN, Pune, Taipei_CWB	Osaka, XiangHe, Pokhara
Urban/Industry	Brookhaven, Billerica, Belsk, GSFC, Ispra, UMBC, Lille, Mexico_City, Moldova, MD_Science_Center, Wallops	Athens_Noa, Shirahama, Leipzig
Biomass Burning	Abracos_Hill, Alta_Floresta, Cuiaba, Concepcion, Los_Fieros, Mongu, Senanga, Skukuza, Zambezi	Bonanza_Creek, Etosha_Pan, Lake_Argyle

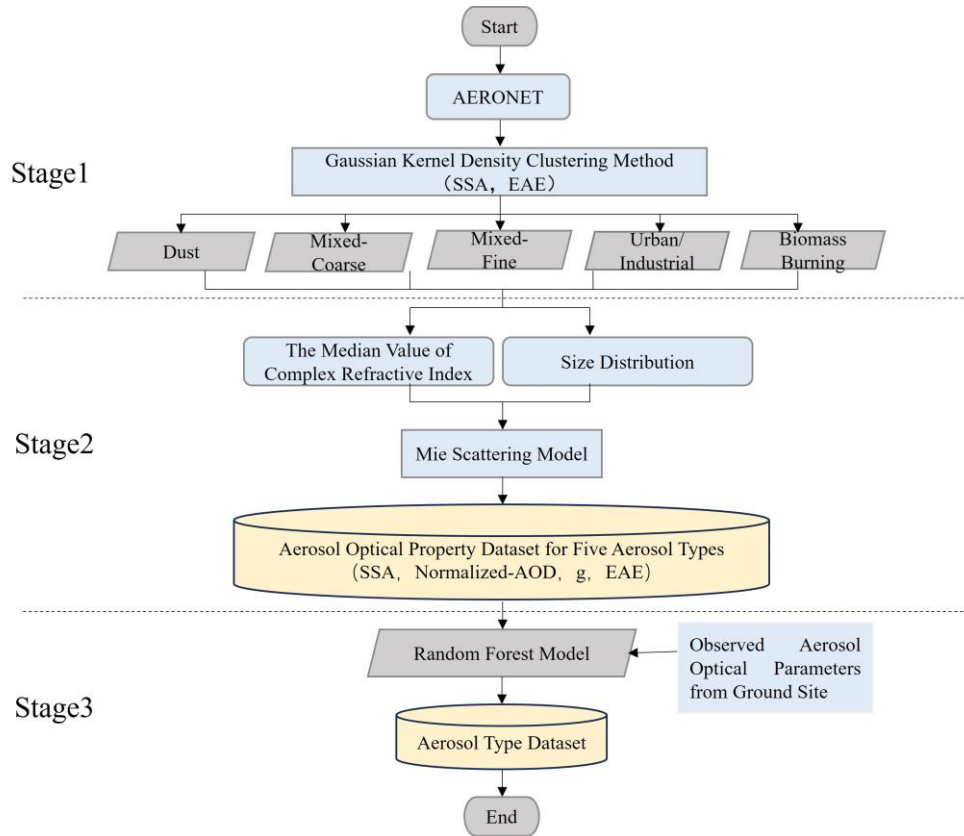
152 **Table 2.** The optical and microphysical properties for aerosol type identification.

	Parameters	Variables (band waves)
Optical Properties	Ångström Exponent (AE)	EAE (440-870) <sup>1</sup>
	Aerosol Optical Depth (AOD)	AOD (440,675,870,1020) <sup>1</sup>
	Single Scattering Albedo (SSA)	SSA (440,675,870,1020) <sup>1</sup>
	Asymmetry Parameter	g (440,675,870,1020) <sup>1</sup>
	Imaginary Part of the Complex Refractive Index	REFI (440,675,870,1020) <sup>1</sup>
Microphysical Properties	Real Part of the Complex Refractive Index	REFR(440,675,870,1020) <sup>1</sup>
	Effective Radius	EffRad-F <sup>2</sup> , EffRad-C <sup>2</sup>
	Standard Deviation of Effective Radius	StaDev-F <sup>2</sup> , StaDev-C <sup>2</sup>
	Size Distribution	Vol-Con (0.05-15µm)

153 Note: <sup>1</sup> refers to wavelength in nm; <sup>2</sup> refers to different modes; EAE is Extinction Ångström Exponent; REFI is Imaginary Part of the  
 154 Complex Refractive Index; REFR is Real Part of the Complex Refractive Index; F refers to fine mode; C refers to coarse mode; EffRad is  
 155 Effective Radius; StaDev is standard deviation; Vol-Con is Volume concentration.

### 156 **3. Methods**

157 A new aerosol classification typing hybrid approach that provides insight into  
158 spatiotemporal variations in aerosol pollution and climate impacts on a global scale is  
159 proposed in this study. In this approach, an aerosol optical properties database using  
160 Mie scattering model was built for calculating rapidly unique aerosol-type features.  
161 Additionally, the approach introduced, for the first time, the median value of the  
162 complex refractive index (CRI) as the criterion for identifying the aerosol type. CRI, a  
163 key microphysical characteristic of aerosols, plays a significant role in determining  
164 their intrinsic optical properties, such as their ability to scatter and absorb light (Raut  
165 and Chazette, 2008). The CRI is also vital for determining aerosols' chemical and  
166 physical compositions (Dubovik and King, 2000) and the CRI value is known for pure  
167 aerosol components (Nandan et al., 2021). Unlike the mean, the median CRI value is  
168 employed in this research for it represents the central tendency of data, especially  
169 beneficial in skewed distributions or when outliers are present. This is particularly  
170 useful when an average value of a specific aerosol-type might be influenced by the  
171 presence of other aerosol types. Moreover, we have selected the aerosol classification  
172 based on the source (as described in Section 1), according to the parameters applied in  
173 this study and the requirements for AOD retrieval. Figure 2 shows the working  
174 flowchart of the new hybrid aerosol-type identification approach, including three  
175 stages: aerosol typing preliminary classification, aerosol optical database generation,  
176 and global aerosol typing identification. The details of these three stages are as  
177 follows.



178

179 **Figure 2.** Flow chart of the new hybrid algorithm in aerosol type identification.

180 **3.1 Aerosol typing preliminary classification (Stage 1)**

181 Stage 1 aimed to solve the problem of obtaining a feature parameter dataset for  
 182 the baseline aerosol type. In previous studies, the Gaussian kernel density clustering  
 183 algorithm showed great potential for distinguishing the optical properties of different  
 184 aerosol types and determining their corresponding thresholds rapidly ( Kalapureddy et  
 185 al. 2009; Pathak et al. 2012). The high concentration value in each cluster generally  
 186 represents the dominant pattern of a specific aerosol type, particularly the data within  
 187 the window, taking the cluster centroid as the center and a specific distance as the  
 188 radius. Preliminary aerosol-type datasets can be generated by digging deep into the  
 189 distribution information of the effective radius, variance, and refractive index of the  
 190 data within the window. The spectral absorbability and particle size of aerosols guide  
 191 the identification of dust, carbonaceous, or hygroscopic aerosols; SSA indicates the  
 192 absorption of aerosol particles; and EAE describes aerosol particle size (Giles et al.,  
 193 2012). Consequently, in this study,  $SSA_{440nm}$  and  $EAE_{440-870nm}$  of 47 AERONET sites



194 and the Gaussian kernel density clustering method was used to estimate the relative  
 195 densities and determine the primary patterns of the dominant aerosol types; here, the  
 196 aerosol type was classified as a dust aerosol. Eqs. (1) and (2) represent the kernel  
 197 density and Gaussian kernel density clustering methods (Rosenblatt, 1956).

$$198 \quad f_{X(v)} = \frac{1}{L} \sum_{i=1}^L k_{\sigma} \left( \frac{\bar{x} - \bar{x}_i}{\sigma} \right), \quad (1)$$

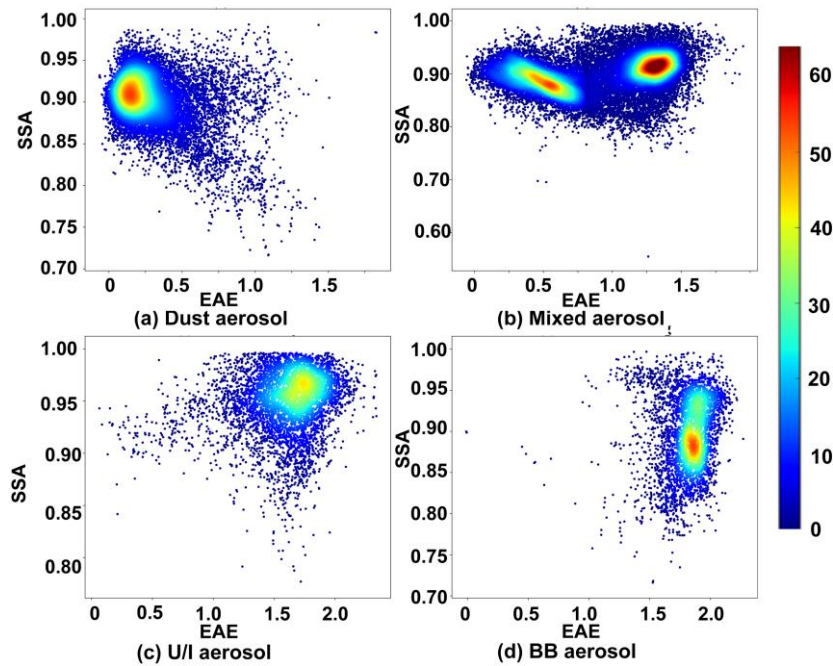
199 where  $f_{X(v)}$  denotes the kernel density and  $k_{\sigma}$  indicates the kernel function.  $x_1,$   
 200  $x_2 \dots x_L$  are the sample points of independent identical distribution. Mathematically,  
 201 kernel functions are symmetric, normalized, and sample-centric when used for density  
 202 estimation; this is best described by the Gaussian kernel equation given by Eq. (2).

$$203 \quad k_{\sigma} = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{|\bar{x} - \bar{x}_i|^2}{2\sigma^2}\right), \quad (2)$$

204 where  $\sigma$  is the kernel size used as a smoothing factor (Moraes et al., 2021).

205 The mixed aerosols comprised fine- and coarse-mode aerosols, indicated by  
 206  $EAE > 0.8$  and  $EAE \leq 0.8$ , respectively. Figure 3 shows the clustering distribution of  
 207 EAE and SSA using the Gaussian kernel density clustering method for different  
 208 aerosol types at the 47 AERONET sites. For the dust aerosol cluster, the density core  
 209 area EAE was 0.1–0.3, and SSA was 0.89–0.94, implying that it contained many  
 210 coarse aerosol particles with moderate absorptivity. Furthermore, the mixed aerosols  
 211 had two distinct centers: one for the coarse-mode aerosols with a median EAE value  
 212 of 0.4, indicating that the cluster contained massive high-absorption aerosols, and the  
 213 other for fine-mode aerosols with a median EAE value of 1.3. Low-absorption  
 214 aerosols were dominant in the cluster, similar to U/I aerosols. Additionally, the density  
 215 core region EAE of U/I aerosol was 1.5–1.8, and SSA was 0.94–0.97, implying the  
 216 dominance of fine and low-absorption aerosols. Conversely, BB aerosols had two  
 217 indistinct centers. This is because, during biomass combustion, gas and particulate  
 218 matter emissions are limited by the combustion conditions, divided into combustion  
 219 and simmering. Combustion produces black smoke, and simmering produces white  
 220 smoke. Combustion, such as burning flames (grass) with high black carbon content,

221 has a strong absorption capacity, resulting in a low SSA. Simmering, such as burning  
222 wood (i.e., trees), tends to be smoldering, lasts longer, has a weaker absorption  
223 capacity, and has a higher SSA value. Therefore, despite possessing different  
224 absorption characteristics, BB aerosols are defined as one aerosol type with an  
225 unseparated center of combustion and simmering.



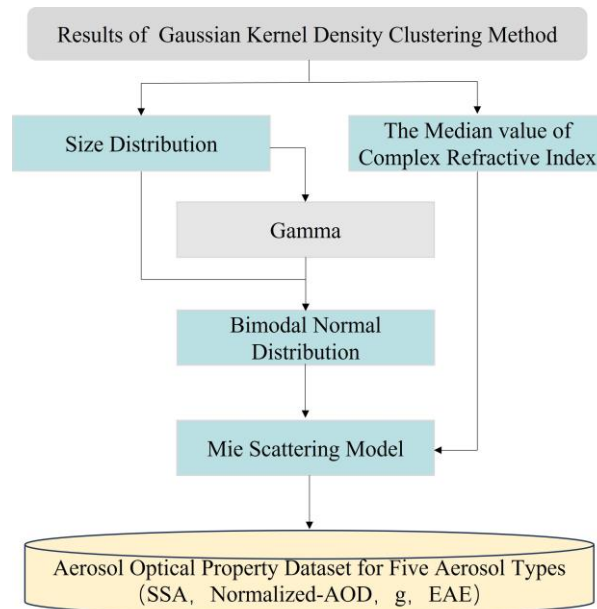
226

227 **Figure 3.** The clustering distribution of EAE and SSA using the Gaussian kernel density clustering  
228 method for different aerosol types.

### 229 3.2 Aerosol optical database generation (Stage 2)

230 In stage 2, the aerosol optical parameter database was built using the aerosol size  
231 distribution parameters, CRI, and Mie scattering model. The main reasons for  
232 constructing an aerosol optical parameter database instead of using the AERONET  
233 data directly are as follows: 1) many data are missed in AERONET, particularly those  
234 for sites dominated by biomass combustion, which does not meet the requirements of  
235 machine learning methods or traditional aerosol type identification algorithms; 2)  
236 Calculating the optical properties of aerosols based on a fixed refractive index can  
237 accurately determine aerosol types. Therefore, once the aerosol spectral distribution  
238 parameters, such as effective radius, variance, and refractive index, are determined in  
239 stage 1, the aerosol optical parameter database can be constructed using the Mie

240 scattering model in stage 2, assuming that aerosols are spherical particles. The Mie  
 241 scattering model, known for its simplicity and practicality, provides an analytic  
 242 solution to Maxwell's equations for light scattering by ideal spherical particles. It  
 243 efficiently depicts the scattering and absorption properties of aerosols in the  
 244 atmosphere, serving as fundamental basis of radiative transfer, Lidar, and optical  
 245 particle characterization (Ma et al.,2007; Bian et al., 2017; Michael et al., 1994).  
 246 Figure 4 details the creation of aerosol optical database, featuring four major  
 247 parameters (normalized-AOD, EAE, SSA, and  $g$ ) at four wavelengths (440, 675, 870,  
 248 and 1020 nm, respectively).



249

250 **Figure 4.** The diagram of building aerosol optical property database.

251 **Table 3.** Size distribution parameters of five aerosol types in coarse and fine mode (unit:  $\mu\text{m}$ )

Aerosol type	REFF-fine	REFF-coarse	Std-fine	Std-coarse
Dust	0.05-0.42	1.3-2.65	0.5-0.8	0.4-0.7
Mixed-coarse	0.05-0.25	1.25-3.5	0.4-0.8	0.4-0.7
Mixed-fine	0.05-0.27	1.2-4.5	0.3-0.6	0.5-0.8
U/I	0.05-0.26	1.45-3.5	0.3-0.6	0.5-0.8
BB	0.05-0.17	1.35-4.5	0.3-0.5	0.5-0.8

252

253

254

255

Table 3 presents the aerosol size distribution parameters, including the effective radius and standard deviation range for the five aerosol types in coarse and fine modes, which were derived from the data window set by the Gaussian kernel density clustering algorithm. These aerosol size distribution parameters and the

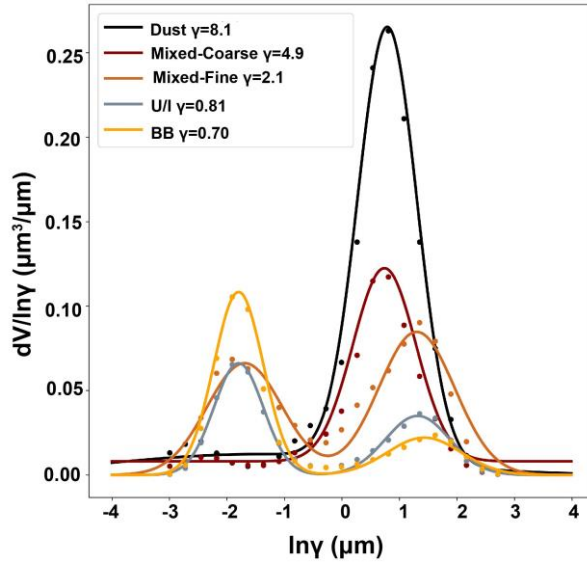
256 median CRI value were utilized to construct the optical database for the Mie  
 257 scattering model. Many studies proven it is a reliable model with the advantage of  
 258 lower computing load and high calculation accuracy (Zhao et al., 2008; Fu et al., 2009;  
 259 Quirantes et al, 2019; Nandan et al., 2021).

260 The Mie scattering model has various size distribution functions, including log-  
 261 normal, power-law, and bimodal log-normal distributions, which describe the aerosol  
 262 type. According to the particle radii provided by AERONET, the size distributions of  
 263 different aerosol types can be divided into coarse and fine modes. The bimodal log-  
 264 normal function [Eq. (3)] is reportedly the most suitable size distribution function for  
 265 modeling aerosol particle size distribution (Remer et al., 2009):

$$266 \quad n(r) = \text{constant} \times r^{-4} \left\{ \exp\left(-\frac{(\ln r - \ln r_{g1})^2}{2 \ln^2 \sigma_{g1}}\right) + \gamma \exp\left(-\frac{(\ln r - \ln r_{g2})^2}{2 \ln^2 \sigma_{g2}}\right) \right\}, \quad (3)$$

267 where  $n(r)$  represents particle count at various radii; constant is obtained by fitting;  
 268 While  $r_{g1}$  and  $r_{g2}$  denote the radii,  $\sigma_{g1}$ , and  $\sigma_{g2}$  are variances for coarse and fine aerosol  
 269 modes, respectively;  $\gamma$ , defined by volume distribution, represents the coarse-to-fine  
 270 mode ratio in bimodal normal distribution model, fitted using AERONET's volume  
 271 distribution data, which averages standard aerosols post-clustering at training sites.

272 Figure 5 shows the volume distributions of five aerosol types, showing dust  
 273 aerosols with a peak  $\gamma$  of 8.1 and radii concentrated around 1.5–2.0  $\mu\text{m}$ . Additionally,  
 274 the mixed-coarse aerosol with a radius in the range of 0.04–0.2  $\mu\text{m}$  and 4.9 as the  
 275 maximum value of  $\gamma$ . The mixed-fine aerosol had two obvious peaks: one with a large  
 276 radius, namely the coarse mode, with a radius of 2.2–3  $\mu\text{m}$  and 2.1 as the peak point  
 277 of  $\gamma$ ; a second with a small radius of 0.1–0.22  $\mu\text{m}$  and 0.14 as the peak point of  $\gamma$ .  
 278 Moreover, the volume distributions of U/I and BB aerosols were similar. Both had a  
 279 relatively low range of  $\gamma$  values at large radii and relatively high values at small radii,  
 280 with peak values of 0.81 and 0.7 for U/I and BB aerosols, respectively.



281

282 **Figure 5.** Volume distribution of five aerosol types.

283 **Table 4.** Real and imaginary index of CRI for five aerosol types (Bands:440/675/870/1020 nm).

Aerosol Type	Imaginary Index	Real Index
Dust	0.003396/0.000731/0.000639/0.000597	1.4584/1.4681/1.4513/1.4376
Mixed-coarse	0.005766/0.002921/0.002383/0.002043	1.4291/1.4787/1.4745/1.4695
Mixed-fine	0.01075/0.008444/0.009147/0.008955	1.5001/1.5044/1.5056/1.4977
U/I	0.004315/0.004331/0.004419/0.004432	1.4372/1.4280/1.4264/1.4214
BB	0.01828/0.017862/0.018125/0.017858	1.5051/1.5190/1.5228/1.5185

284 The CRI is an inherent optical property of aerosols. Aerosols in the real  
 285 atmosphere are usually mixed with different types of particles, which a single  
 286 refractive index cannot identify; however, the CRI represents the entire aerosol model  
 287 in the atmosphere (Redemann et al., 2000). Ideally, the CRI and aerosol components  
 288 can be mutually determined (Wu et al., 2021). The CRI can effectively characterize  
 289 the main properties of the aerosols and accurately quantify the difference between  
 290 aerosol-type identification algorithms. Table 4 depicts the CRI standard values for the  
 291 five aerosol types obtained by calculating the median value of the CRI of the  
 292 dominant aerosol type after Gaussian kernel density clustering. These values were  
 293 used as a baseline for identifying the aerosol types in subsequent studies. As presented  
 294 in Table 4, the minimum imaginary index part is represented by the dust aerosol with  
 295 CRI of 0.003396, 0.000731, 0.000639, and 0.000597 at 440, 675, 870, and 1020 nm,

296 respectively, owing to the weakest absorption of dust aerosols. Moreover, the  
 297 imaginary index part of the mixed-fine aerosols (0.01) was close to that of the BB  
 298 aerosols (0.02) because of their similar absorption properties.

299 Lastly, by fixing the CRI, changing the size distribution, and using the Mie  
 300 scattering model, we generated the aerosol optical property database for five aerosols,  
 301 including the data for normalized-AOD, EAE, SSA, and  $g$ . In the aerosol optical  
 302 property database, normalized AOD is the value obtained after eliminating the  
 303 influence of the aerosol concentration. The AOD was obtained from the extinction  
 304 cross section ( $C_{ext}$ ) calculated using the Mie scattering model in Eqs. (3) and (4),  
 305 where  $\beta_{ext}$  is the extinction coefficient,  $n(r)$  is the aerosol spectral distribution, and  
 306  $N(z)$  is the variation of aerosol concentration with height. Notably, the effect of  
 307 aerosol concentration needs to be removed from the AOD when referring to aerosol  
 308 optical properties. The AOD was normalized by dividing the aerosol optical thickness  
 309 at the four wavelengths by the optical thickness at 440 nm. The other parameters  
 310 (EAE, SSA, and  $g$ ) were calculated using Eqs. (6) – (8).

$$311 \quad \beta_{e/s} = \int_{\gamma_{min}}^{\gamma_{max}} C_{ext/sca} n(r) dr , \quad (4)$$

$$312 \quad \tau_{e/s} = \int_0^{Z_{top}} \beta_{ext/sca} N(z) dz, \quad (5)$$

$$313 \quad EAE_{440-870nm} = -\frac{\ln(\tau_{440nm}) - \ln(\tau_{870nm})}{\ln(440) - \ln(870)} , \quad (6)$$

$$314 \quad SSA = \frac{\tau_s}{\tau_e} , \quad (7)$$

315 and

$$316 \quad g = \langle \cos\Theta \rangle = \frac{1}{2} \int_{-1}^1 p(\cos\Theta) \cos\Theta d \cos\Theta , \quad (8)$$

317 where  $\tau_{440}$  and  $\tau_{870}$  are the extinction optical depths of the aerosol at 440 and 870 nm,  
 318 respectively,  $EAE_{440-870}$  nm is the extinction Ångström index from the 440 to 870 nm  
 319 band, and  $\Theta$  denotes the scattering angle.

### 320 **3.3 Global aerosol type identification and validation (Stage 3)**

321 In stage 3, the random forest model was introduced to the aerosol-type  
322 identification algorithm. The random forest model is an integrated model based on  
323 classification and regression trees, in which multiple trees are aggregated using  
324 majority voting and averaging for classification and regression (Breiman, 2001). The  
325 model has a high prediction accuracy, excellent tolerance for abnormal values and  
326 noise, and a hard overfit. In a comparison by Fernandez (2014), the random forest  
327 algorithm ranked as the top performer among 179 classification algorithms. In  
328 addition, the evaluation matrix was brought into this study, and it further  
329 quantitatively assesses the performance of the Gaussian density clustering algorithm  
330 and the new hybrid algorithm. The metric indexes include accuracy, recall, precision,  
331 and F-scores (Reddy et al., 2022). Here, the indexes are adjusted to micro-precision,  
332 micro-recall, micro-F1-score, and accuracy to solve the multi-classification problem.  
333 Micro refers to the weighted average of the five aerosol types rather than the  
334 arithmetic mean, due to the large difference in sample size among the five aerosol  
335 types, the arithmetic mean is highly susceptible to the influence of very large or very  
336 few sample size aerosol types.

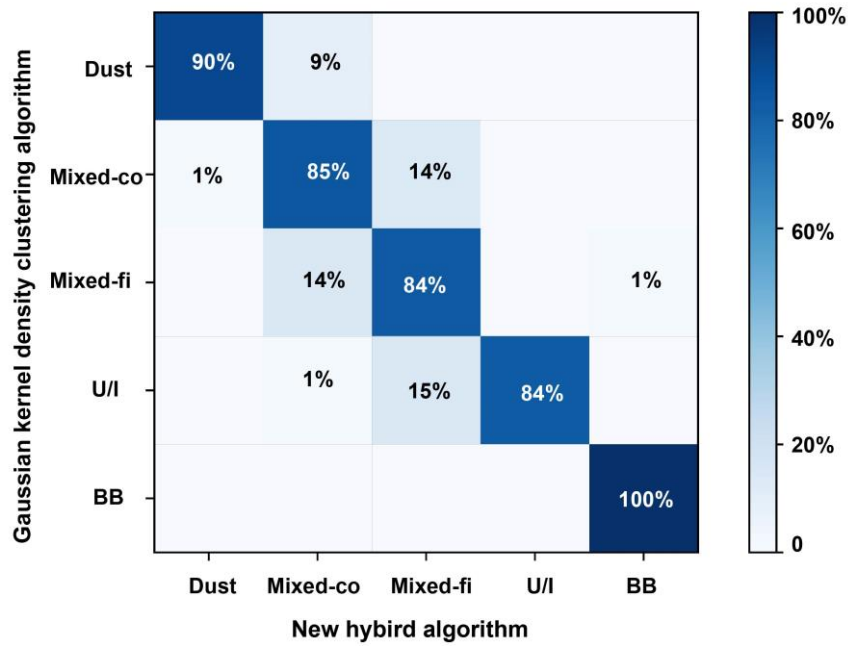
337 The input parameters for random forest model training, including  $SSA_{440nm}$ ,  
338  $SSA_{675nm}$ ,  $SSA_{870nm}$ ,  $SSA_{1020nm}$ ,  $g_{440nm}$ ,  $g_{675nm}$ ,  $g_{870nm}$ ,  $g_{1020nm}$ , normalized  $AOD_{675nm}$ ,  
339  $AOD_{870nm}$ ,  $AOD_{1020nm}$ , and  $EAE_{440-870nm}$ , were selected from the aerosol optical  
340 property database, and the expected output values were the specific aerosol types. The  
341 random forest model was optimized, and the parameters were determined using the  
342 grid-searching method. The parameters, including `n_estimators` (classifier),  
343 `max_features` (maximum feature value), and `min_samples_leaf` (minimum number of  
344 samples for nodes), were set as 160, 10, 12, and 12, respectively. Then, based on the  
345 trained and optimized model, aerosol typing of any AERONET site in different  
346 regions of the world can be identified quickly. Generating the aerosol-type  
347 distribution map on a global scale is vital for regional and global climate studies as  
348 well as ground remote sensing.

## 349 **4 Results**

### 350 **4.1 Algorithm comparison**

351 To demonstrate the effectiveness of the new hybrid algorithm, its performance  
352 was compared with that of the Gaussian kernel density clustering algorithm. Figure 6  
353 shows the confusion matrix between the new hybrid and Gaussian kernel density  
354 clustering algorithms in identifying aerosol types. The results of the new hybrid  
355 algorithm showed 90% consistency with that from the Gaussian kernel density  
356 clustering algorithm, in delineating dusty aerosols, indicating that its efficiency in  
357 identifying dust. For mixed-coarse aerosols, the consistency reached 85%, with 14%  
358 identified as mixed-fine aerosols, 1% as dust by the new hybrid algorithm, and 15%  
359 as mixed-coarse aerosols by the Gaussian kernel density clustering algorithm.  
360 Similarly, for mixed-fine aerosols, both algorithms showed 84% consistency, with 14%  
361 identified as a mixed-coarse aerosol by the new hybrid algorithm and as a mixed-fine  
362 aerosol by the Gaussian kernel density cluster algorithm. Furthermore, both  
363 algorithms identified 84% of U/I aerosols correctly, with the remaining 16% identified  
364 as mixed aerosols (fine and coarse). Lastly, the classification of BB aerosols using  
365 these two methods was the same. Overall, the Gaussian kernel density clustering and  
366 new hybrid algorithms were consistent in dust, mixed-coarse, U/I, and BB aerosol  
367 identification.





368

369 **Figure 6.** The confusion matrix between Gaussian kernel density clustering and new hybrid  
 370 algorithm.

371 Table 5 shows the metric index value of the random forest algorithm in the new  
 372 hybrid algorithm. The micro-precision, micro-recall, micro-F1 score, and accuracy are  
 373 0.95, 0.89, 0.91, and 0.89, respectively. These metrics are derived from the core  
 374 values of the window, as determined by the Gaussian density clustering algorithm.  
 375 Consequently, the strong performance of these indicators further confirms the efficacy  
 376 and reliability of the newly developed hybrid algorithm.

377 **Table 5.** Matrix evaluation between new hybrid classification algorithm and Gaussian kernel  
 378 density clustering algorithm

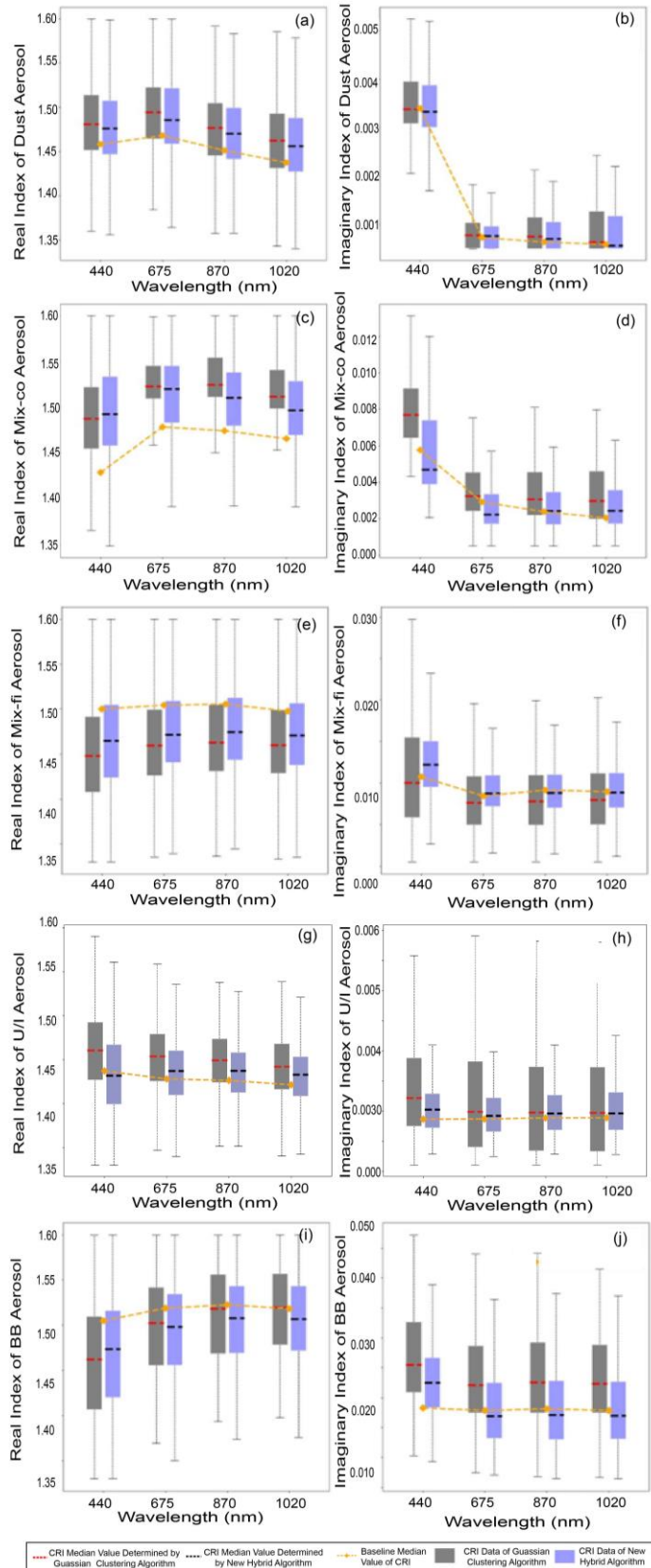
	Micro-Precision	Micro-Recall	Micro-F1-Score	Accuracy
New Hybrid algorithm	0.95	0.89	0.91	0.89

379 As described in the Methods section, a specific aerosol type theoretically has a  
 380 fixed CRI owing to its constant composition. The CRI characterizes the mixture  
 381 composition of aerosol particles and is a key parameter controlling the inherent  
 382 scattering and absorption characteristics of aerosol particles. To further analyze the  
 383 accuracy of the new algorithm, the aerosol CRI was applied as a key criterion for  
 384 aerosol identification. The CRI has two parts: imaginary and real. The imaginary part  
 385 indicates radiation absorption by aerosols, with a small value signifying a small

386 absorption. Because the radiation of aerosols is more dependent on the imaginary than  
387 the real part, the imaginary part is essential for inferring the optical properties and  
388 aerosol types. Hence, we compared the real and imaginary parts of the CRI calculated  
389 using the new hybrid and Gaussian kernel density clustering algorithms.

390 Figure 7 shows box plots of the aerosol CRI for dust, mixed-coarse, mixed-fine,  
391 U/I, and BB aerosols using the new hybrid classification and Gaussian kernel density  
392 clustering algorithms. Based on the principle that the CRI of aerosols is fixed under  
393 ideal conditions, the closer the median value of the CRI of the identified aerosol type  
394 is to the median value of the benchmark CRI, the more accurate the identification  
395 method. As shown in Figures 7 (a) and (f), the median values of the CRI real part for  
396 dust aerosol are in the range 1.45–1.53 at four bands, and those of the imaginary part  
397 are 0.003–0.004 at 440 nm; further, the values in other bands decrease rapidly as  
398 wavelength increases. The imaginary part of CRI represents the absorption of light by  
399 the aerosol, with a small absorption indicating strong scattering. The results of the  
400 imaginary part are consistent with the spectral dependence properties of dust-based  
401 aerosols according to the wavelength. This is primarily because dust aerosols,  
402 composed of clay, quartz, and hematite, exhibit strong absorption in the blue band  
403 (440 nm) and low absorption in the visible and near-infrared bands. For the dust  
404 aerosols, the CRI determined by the two methods did not differ much. However, the  
405 median value of the CRI obtained using the new hybrid algorithm was slightly closer  
406 to the benchmark CRI than that obtained using the Gaussian kernel density clustering  
407 algorithm for dust aerosols. Therefore, the new hybrid algorithm was concluded to be  
408 more accurate in identifying dust aerosol.

409 Figures 7 (b) and (g) show the median values of the CRI real part for mixed-  
410 coarse aerosol is 1.47–1.55 at four bands using the new hybrid algorithm, but the  
411 imaginary part is 0.004–0.009 at 440 nm. However, the real part is 1.44–1.50 at four  
412 bands determined by the Gaussian kernel density clustering algorithm, and the  
413 imaginary part is 0.006–0.009 at 440nm. The median value of the hybrid algorithm  
414 was closer to the baseline median value than that of the Gaussian kernel density  
415 clustering algorithm for both the real and imaginary parts.



416

417 **Figure 7.** Box plots of the real index (left) and the imaginary (right) index of the CRI for (a-b)  
 418 dust, (c-d) mixed-coarse, (e-f) mixed-fine aerosol, (g-h) U/I, and (i-j) BB aerosol identified by the  
 419 Gaussian kernel density clustering algorithm and new hybrid algorithm, respectively.

420 Figures 7 (c) and (h) show the median value of the CRI real part for mixed-fine  
421 aerosols determined using the new hybrid and Gaussian kernel density clustering  
422 algorithms, which was 1.42–1.51 at four bands. This result is close to the range (1.44–  
423 1.52) reported by Wu (2021) in Beijing using a random forest algorithm. The median  
424 CRI of the real part at four bands and the imaginary part at the (675-870-1020 nm)  
425 bands were close to the baseline median value for the new algorithm. Additionally, the  
426 median value of the imaginary part was lower than that of the new hybrid algorithm  
427 and further from baseline data for the identifying aerosol type results mixed with 14%  
428 coarse aerosols. Mixed coarse aerosols result in weaker absorption. Hence, the new  
429 hybrid algorithm performed better at identifying mixed-fine aerosols than the  
430 Gaussian kernel density clustering algorithm.

431 Similarly, as seen in Figures 7 (d) and (i), the median value of the CRI real part  
432 for U/I aerosol identified using the new hybrid algorithm was 1.39–1.47. This median  
433 value was lower than that of the mixed-fine aerosols. This is because the real part  
434 indicates the absorption ability of aerosols, and the absorption ability of U/I aerosols  
435 was less than that of mixed-fine aerosols. For the imaginary part also, the new hybrid  
436 algorithm performed slightly better than the Gaussian kernel density clustering  
437 algorithm at the four bands.

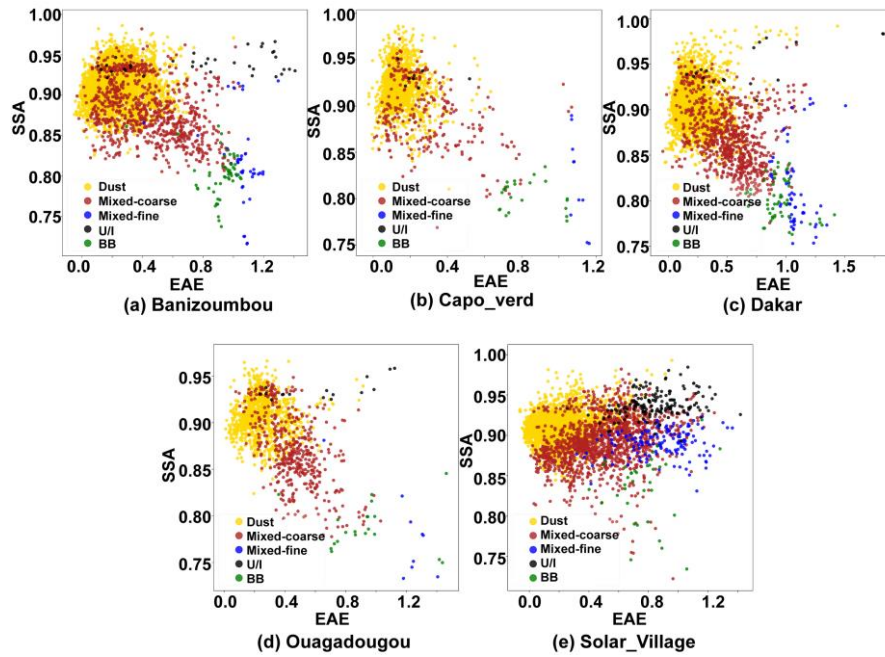
438 For BB aerosols, the median value of the real part generated using the new hybrid  
439 algorithm differed slightly from that generated by the Gaussian kernel density  
440 clustering algorithm. Additionally, the median obtained using the Gaussian kernel  
441 density clustering algorithm was closer to the baseline. Furthermore, when analyzing  
442 the imaginary part, the new hybrid algorithm performed much better than the  
443 Gaussian kernel density clustering algorithm. Even with a 100% concordance rate  
444 between the new hybrid and Gaussian kernel density clustering algorithms in  
445 identifying BB aerosols, the refractive index still differed. This result indicates that 1%  
446 of mixed-fine aerosols classified using the Gaussian kernel density clustering  
447 algorithm were correctly identified as BB aerosols by the new algorithm. Overall,  
448 these results demonstrate that the new algorithm is reliable.

449           Additionally, in this study, the number of 326400 data points from optical  
450 parameters database and 98000 observed data for calculation spans from Jan.1st,1993  
451 to Dec.31st,2021, passing through Gaussian kernel density clustering algorithm and  
452 new hybrid algorithm Python progresses, which is archived on the personal Windows  
453 system computer (Intel® Core™ i7-10710U,16G DDR4 2666MHz, 512G PCIE SSD).  
454 The computational time for the two algorithms indicates the new hybrid algorithm  
455 runs faster than the Gaussian kernel density clustering algorithm with huge quantities  
456 of data and trained in advance, which can obtain aerosol type in 20 seconds, in  
457 contrast, it will take 30 to 40 seconds to obtain aerosol type in one site by using the  
458 Gaussian algorithm.

## 459   **4.2 Aerosol type determination for typical sites**

### 460   **4.2.1 Dust aerosol**

461           Figure 8 shows the aerosol types obtained using the new hybrid algorithm for the  
462 five sites selected for dust aerosol identification. According to the prediction by the  
463 new hybrid algorithm, the aerosols at these five sites mainly contained dust aerosols  
464 along with a small amount of U/I, mixed-fine, and BB aerosols, and a large amount of  
465 mixed coarse aerosols. This shows that other types of aerosols invaded these areas  
466 besides dust aerosol. BB aerosols may have been transferred from the southern  
467 African savannah. Additionally, U/I aerosols could be from industrial cities, such as  
468 Dakar, Abidjan, and Lagos, which are dominated by anthropogenic aerosols and are  
469 close to the AERONET sites.



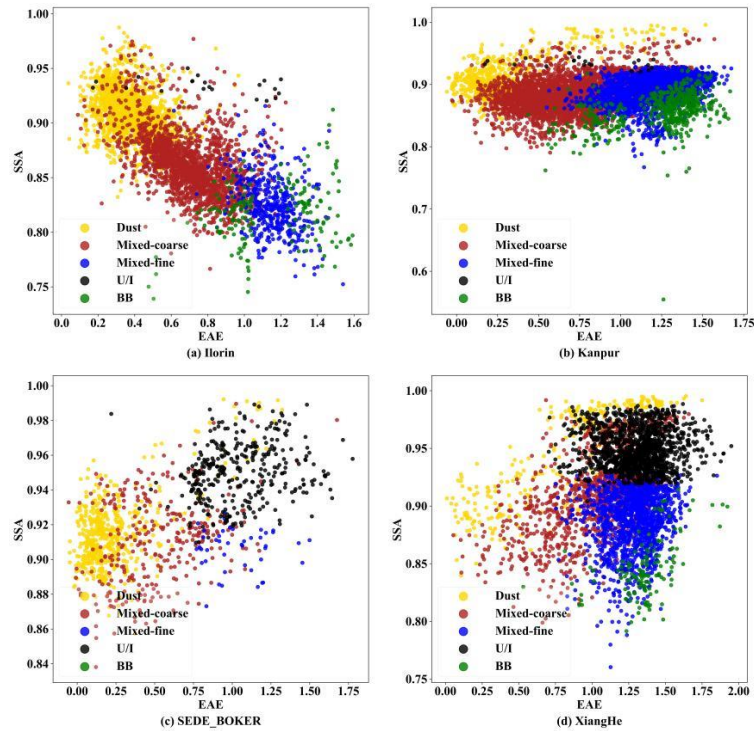
470

471 **Figure 8.** Identification of dust aerosol at dominant aerosol sites.

#### 472 4.2.2 Mixed aerosol

473 Besides Ilorin in Africa, the mixed aerosol AERONET sites, including Kanpur,  
 474 Sede\_Boker, and XiangHe, are in Asia. The aerosol types at these four sites were  
 475 determined using the new hybrid algorithm (Figure 9). Mixed coarse aerosols  
 476 dominated the Kanpur, Ilorin, and Sede\_Boker sites, and mixed fine aerosols  
 477 dominated XiangHe. Part of the dust in Xianghe could be due to the Takla Desert in  
 478 spring and the westerly winds prevailing in western China, which transported dust  
 479 aerosols over long distances. Additionally, the U/I aerosol in Xianghe could be a result  
 480 of human activities, construction emissions, and fuel burning in winter. The BB  
 481 aerosol was traced to the burning of a small amount of biomass in Xianghe, located in  
 482 a suburban area.

483 Furthermore, excluding dust aerosols, we observed BB and U/I aerosols in the  
 484 Kanpur site in the Ganges Basin of India. A certain amount of U/I and dust aerosols  
 485 were also observed in Sede\_Boker, located in the industrial center of Israel, possibly  
 486 from the Arabian desert. Lastly, Ilorin had the most dust and least BB aerosols  
 487 because it is located in central Africa, often affected by the Saharan Desert and  
 488 African savannah.



489  
 490 **Figure 9.** Same as Figure 8 but for Mixed aerosol.

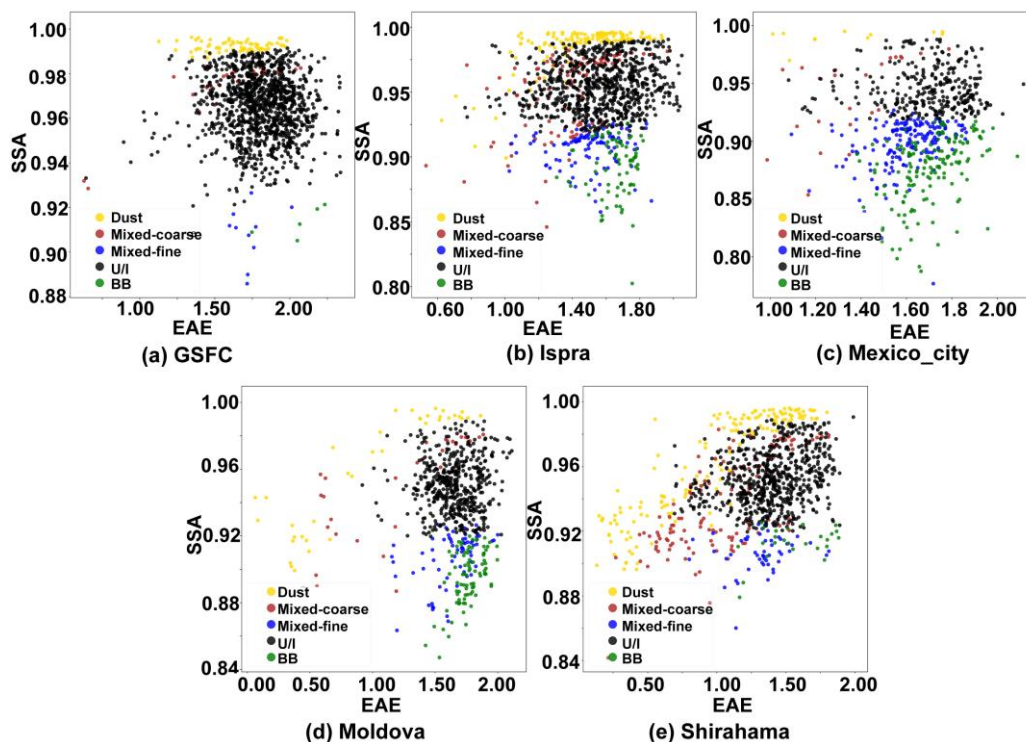
491 **4.2.3 Urban/industrial aerosol**

492 All the selected AERONET sites for evaluating the performance of the new  
 493 hybrid algorithm in terms of U/I aerosol identification are in Europe or North America  
 494 (Figure 10). GSFC is located in the densely populated and industrially developed area  
 495 of Washington in the United States, explaining its complex aerosol type dominated by  
 496 the U/I aerosol followed by a few mixed and BB aerosols and a small amount of dust  
 497 aerosols.

498 Ispra is in Turin, one of Italy’s largest industrial centers. However, dust-type  
 499 aerosols were identified, possibly transported from the Libyan desert when Italian  
 500 winters were controlled by southwesterly winds. Moreover, Mexico, where the  
 501 Mexico City site is located, is an industrialized country with modern industries and  
 502 agriculture, abundant oil production, and a dense population. Nevertheless, we  
 503 identified dust, mixed coarse, and BB aerosols in this site using the new hybrid  
 504 algorithm. These aerosol types could be from the Chihuahuan Desert, an inland desert  
 505 covering 12% of Mexico's area and a major source of coarse and dust aerosols.  
 506 Additionally, the literature shows that Mexico City is surrounded by forested



507 mountains, which experience many wildfires during the dry period between  
 508 November and May; this accounts for BB aerosols in Mexico City (Yokelson et al.  
 509 2007). Finally, the BB aerosols identified at the Moldova site could be attributed to its  
 510 rich vegetation cover.



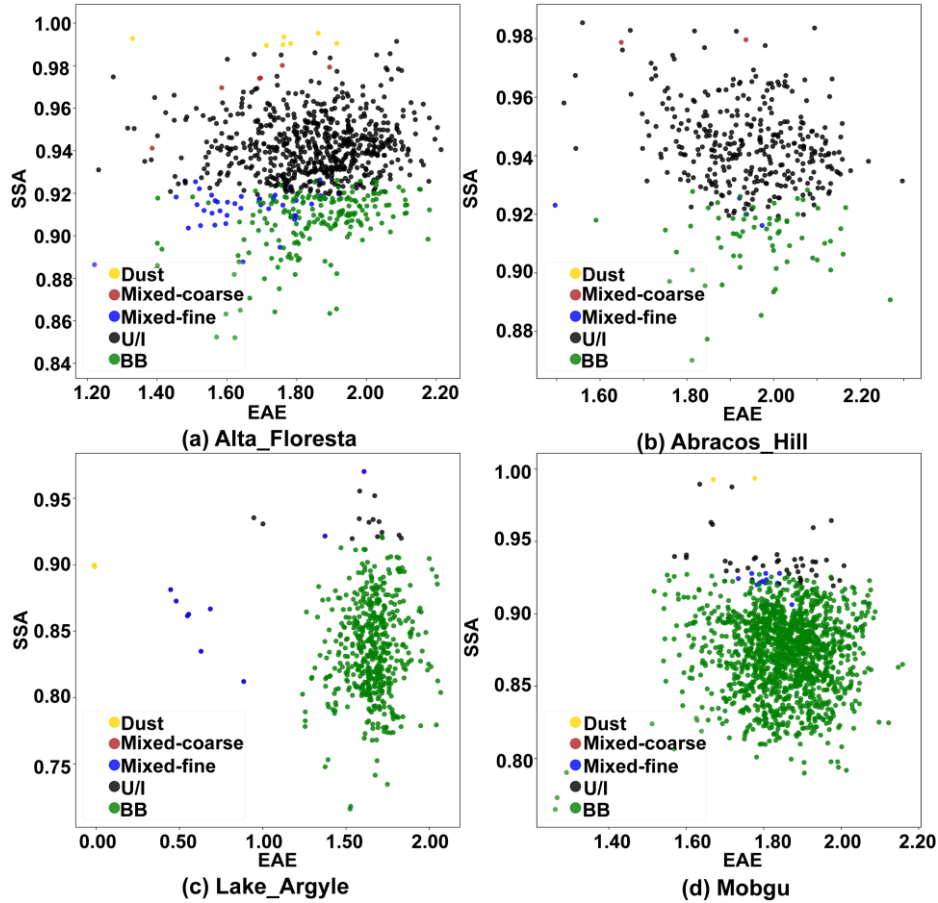
511  
 512 **Figure 10.** Same as Figure 9 but for urban/industrial aerosol.

#### 513 4.2.4 Biomass burning aerosol

514 The selected sites were mainly located in the mountains and highlands. Figure 11  
 515 shows the aerosol types identified using the new hybrid algorithm. Large amounts of  
 516 BB aerosols were identified at all sites. Additionally, a small amount of dust and  
 517 mixed-coarse aerosols were identified at the Alta\_Floresta site, transported over a  
 518 long distance from the Patagonian Desert in Argentina, in southern South America.  
 519 Moreover, the city where the site is located is industrially developed and has a large  
 520 population; therefore, more U/I aerosols were identified using the new hybrid  
 521 algorithm. The geographically close Abracos\_Hill and Alta\_Floresta sites were  
 522 characterized by the same aerosol type and source. Furthermore, one data point in  
 523 Lake Argyle was classified as a dust aerosol. This means that, although the site is  
 524 located on the Kimberley Plateau, Australia has a large desert area, and coarse



525 aerosols still exist. Lastly, a few U/I and several dust-type aerosols were identified at  
 526 the Mongu site, possibly caused by aerosol emissions from nearby cities and dust  
 527 transport from the Saharan Desert.



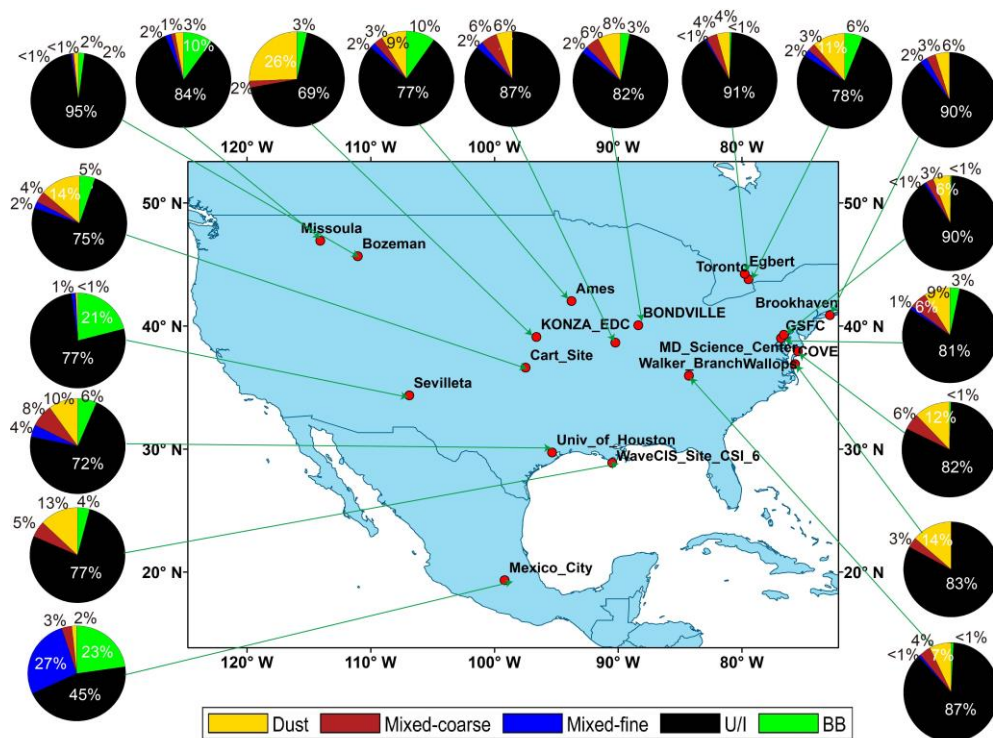
528  
 529 **Figure 11.** Same as Figure 10 but for BB aerosol.

530 **4.3 Aerosol type distribution on a global scale**

531 Given the advantages and accuracy of the new hybrid algorithm in identifying  
 532 aerosol types, we used it to divide the data of AERONET sites in different regions of  
 533 the world to obtain global aerosol type distribution information. The aerosol types of  
 534 each continent are shown in Figures 12-16. Additionally, Figure 17 shows the global  
 535 aerosol-type distribution. Notably, the pie chart was placed on each site in the study,  
 536 which is a "point source" assessment of the aerosol type and does not represent the  
 537 entire region (the size of the pie chart is independent of the optical properties).  
 538 Moreover, the sites were screened, and only those with valid data of > 100 aerosol  
 539 types were considered; however, offshore sites and sites classified as marine aerosol-

540 dominated by other literature were excluded.

541 Figure 12 shows pie charts of the aerosol types for each scanned AERONET site  
 542 in North America. The U/I aerosols, particularly in most mid-eastern regions,  
 543 contained mixed and small amounts of biomass aerosols. Additionally, the AERONET  
 544 sites in large cities, such as Chicago, New York, Toronto, Ottawa, and Boston, had U/I  
 545 aerosols. Many studies have shown that dust aerosols from the Saharan Desert can  
 546 cross the Atlantic Ocean to North America in summer. Moreover, there is an inland  
 547 desert in western North America, the Chihuahua Desert, responsible for a small  
 548 amount of dust and mixed aerosols at the AERONET sites in North America.  
 549 Additionally, wildfires in western North America and household wood burning  
 550 contribute to most BB aerosols yearly. The central region site is affected by the  
 551 environment, with an increased proportion of BB aerosols, and U/I aerosols are still  
 552 prevalent because the site is located in a large city and is densely populated.

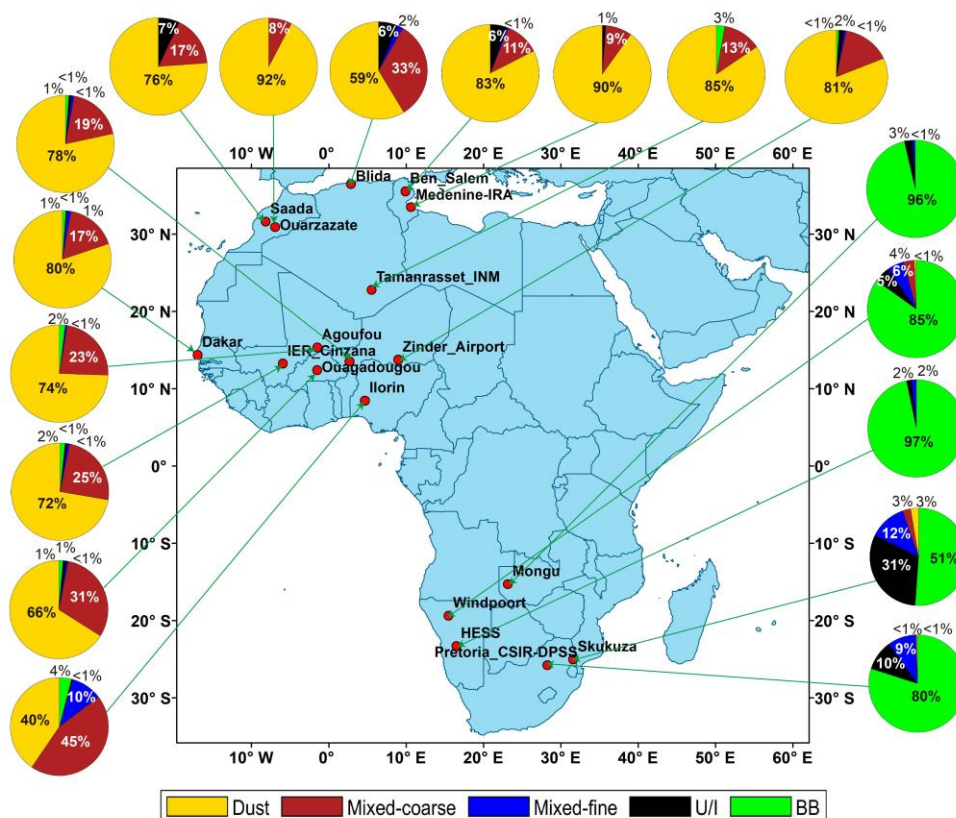


553

554 **Figure 12.** Pie charts of the aerosol types at the major sites of North American.

555 Figure 13 shows the aerosol types in Africa. Northern Africa has the largest desert  
 556 in the world, the Saharan Desert; therefore, dust aerosols dominate north of the  
 557 equator in Africa. However, some AERONET sites in the Sudanese steppe were

558 primarily BB, with some U/I aerosols in nearby urban sites. The Ilorin site is a typical  
 559 mixed aerosol site close to the equator with a small amount of BB aerosols. Most sites  
 560 close to the Atlantic coast were affected by dust aerosols, even those on the islands of  
 561 Capo\_Verde. The reliability of the new model in distinguishing U/I and BB aerosols is  
 562 demonstrated. Sites in Southern Africa, such as Namibia, Botswana, and Zambia, are  
 563 dominated by BB aerosols. Nevertheless, studies have shown the presence of U/I  
 564 aerosols at sites in the urban areas of South Africa. Although U/I and BB aerosols are  
 565 difficult to distinguish, the two can be identified in the context of a large urban  
 566 population and less biomass combustion, thus establishing the model's accuracy.

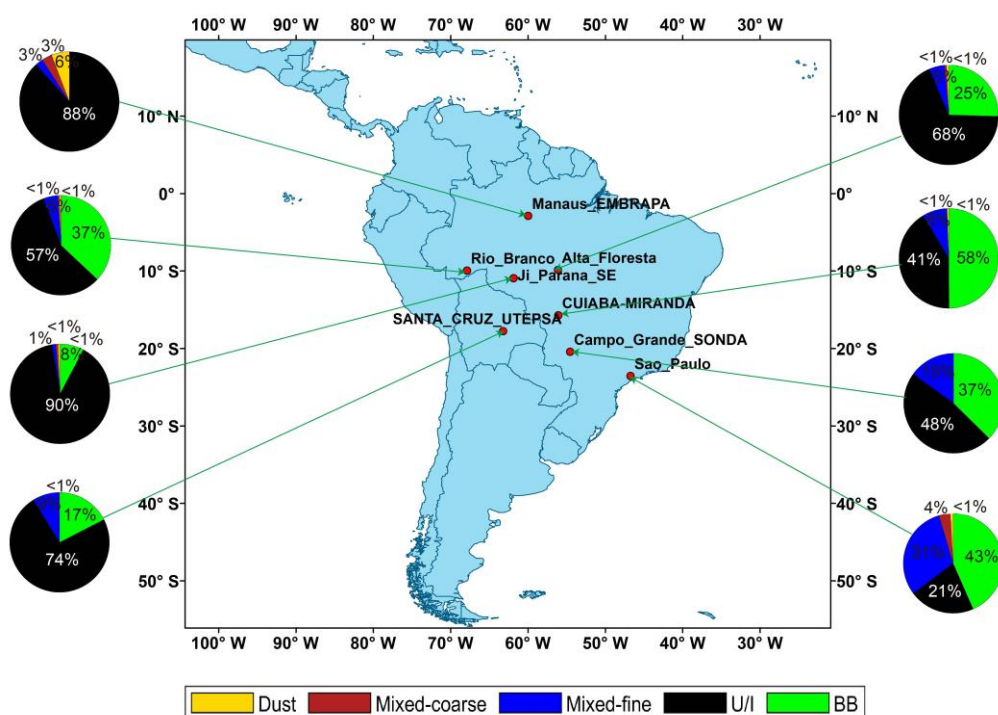


567

568 **Figure 13.** Same as Figure 12 but for Africa.

569 The aerosol types in South America are shown in Figure 14. Here, only eight sites  
 570 met the requirement for valid data >100 aerosol types. South America is mainly  
 571 dominated by mountainous plateaus, and under the influence of the Brazilian warm  
 572 current, many tropical rainforests are distributed in the south; therefore, the  
 573 background aerosols are mainly BB aerosols. As shown in Figure 14, large cities, such  
 574 as Rio Branco, Campo Grande, Manaus, Santa Cruz, and São Paulo, showed an

575 increased proportion of anthropogenic and mixed aerosols because of their large  
 576 population and developed industries. Due to the tropical rainforest climate in southern  
 577 South America, the proportion of BB aerosols increased, such as that at the Cuiaba  
 578 site near the Amazon River. Additionally, the Manaus site contained a small amount of  
 579 dust aerosols that were presumably transported across the Atlantic Ocean from  
 580 African dust at the same latitude.



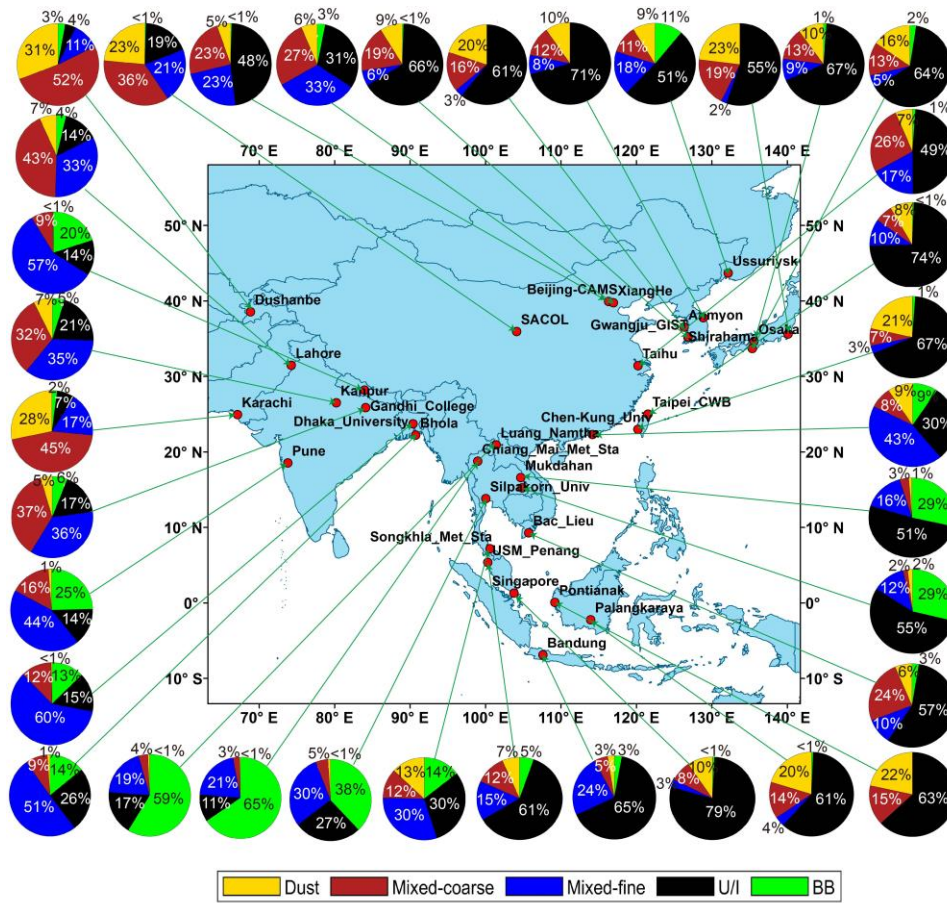
581

582 **Figure 14.** Same as Figure 12 but for South America.

583 The aerosol types in Asia are shown in Figure 15. In western Asia, influenced by  
 584 the Indian Desert, sites on the Indian Peninsula were dominated by coarse-particle  
 585 aerosols, including dust and mixed coarse aerosols. Kanpur and Pune are densely  
 586 populated cities in India, with more mixed-fine aerosols produced by human activities.  
 587 Additionally, in Southeast Asia, all sites contained BB aerosols, consistent with  
 588 Hamill (2014). This is because of the abundance of tropical rainforests in Southeast  
 589 Asia. Moreover, some urban sites, such as Singapore and Penang, had large numbers  
 590 of U/I and mixed-fine aerosols. The coastal areas of East Asia, which are densely  
 591 populated and industrially developed, were mainly dominated by U/I aerosols.



592 Moreover, dust aerosols appeared at these sites due to dust transported from the  
 593 Taklamakan Desert in East Asia.

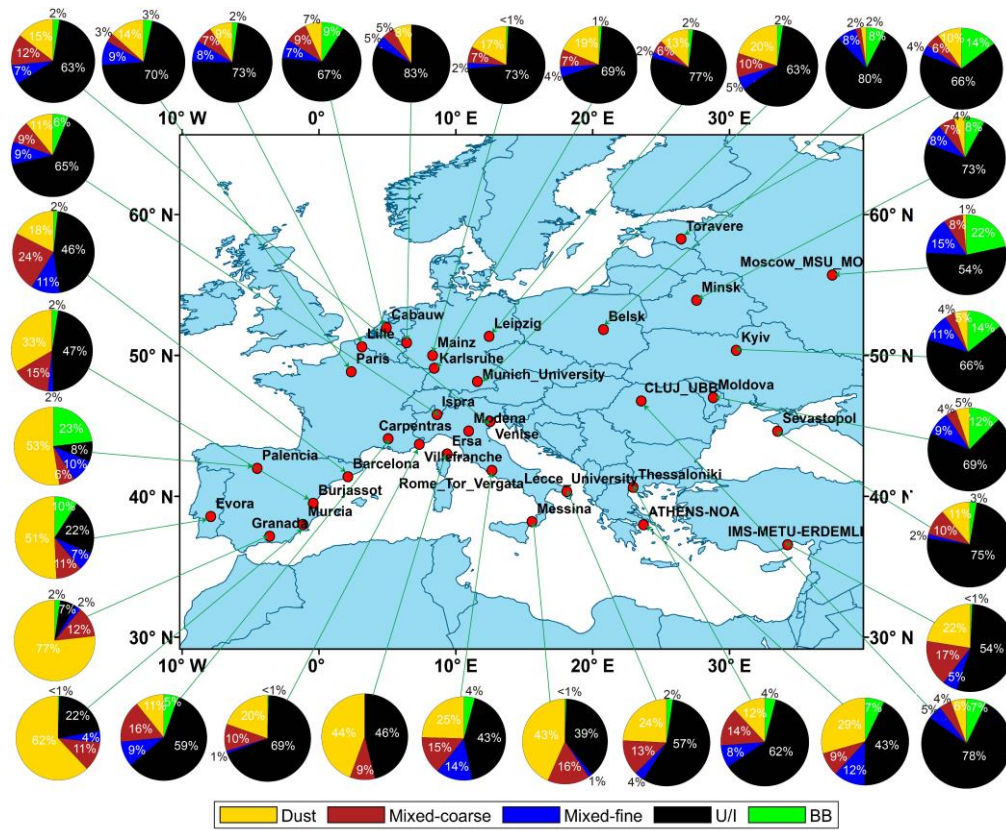


594  
 595 **Figure 15.** Same as Figure 12 but for Asia.

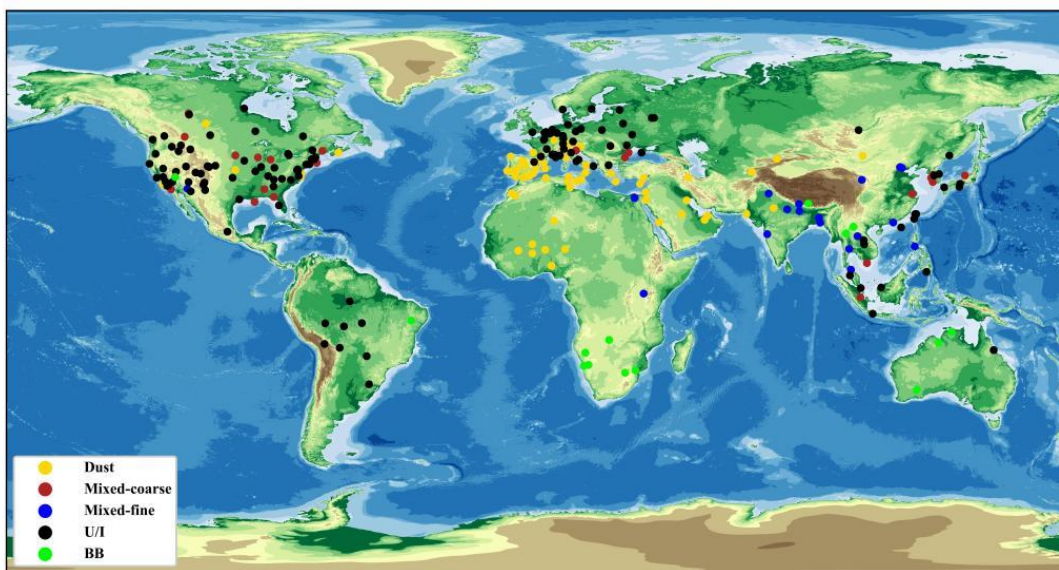
596 The inland areas of East Asia have a smaller population than the coastal areas;  
 597 therefore, the proportion of U/I aerosols was small, and that of mixed aerosols was  
 598 high. Generally, mixed aerosols are more easily overestimated than U/I aerosols;  
 599 however, the new hybrid algorithm identified a larger proportion of U/I aerosols than  
 600 mixed aerosols at Asian sites. Therefore, this new hybrid algorithm can be considered  
 601 for improving the classification of mixed aerosols versus U/I aerosols.

602 Similarly, southern Europe, which is close to the Saharan and Arabian deserts,  
 603 was dominated by dust aerosols, with small amounts of mixed and U/I aerosols.  
 604 Northern European sites have many cities and a large population; therefore, the  
 605 aerosol type was mainly U/I aerosols, identified using the new hybrid algorithm  
 606 (Figure 16). Additionally, small amounts of BB aerosols were identified at most sites  
 607 in Europe because of olive groves in agricultural lands in the EU, which produce 91%

608 of the world's olive oil (Lopez-Pineiro et al., 2011). Papadakis et al. (2015) suggested  
 609 that the biomass produced from olive oil is used for heating and industry, and its  
 610 combustion produces carbonaceous aerosols, considered the major source of fine  
 611 particle aerosols in Europe during winter (Puxbaum et al., 2007).



612  
 613 **Figure 16.** Same as Figure 12 but for Europe.



614  
 615 **Figure 17.** Global dominant aerosol type distribution based on AERONET sites.

616 The global distribution of dominant aerosols in the AERONET site is shown in  
617 Figure 17. The graph does not include marine aerosols. There are more aerosol sites  
618 on the global map than those on each continent because AERONET sites with > 5  
619 years of data were selected for the global map; however, sites with > 100 valid data  
620 points were required for each continent. The global distribution map shows that many  
621 BB aerosols were distributed between 20°N and 20°S. This is because this region has  
622 a predominantly tropical rainforest climate, with many tropical rainforests and more  
623 carbon-containing aerosol emissions. This finding is consistent with those from  
624 previous studies that found that global BB aerosols mainly originate from Africa  
625 (approximately 52%), followed by South America (approximately 15%), equatorial  
626 Asia (approximately 10%), boreal forests (approximately 9%), and Australia  
627 (approximately 7%) (Van G. R. et al., 2010). Furthermore, the global distribution map  
628 shows a clear distribution band of dust aerosols between 5°N and 35°N, originating  
629 from the Saharan Desert in Africa and the Saudi Arabian Desert in Western Asia,  
630 which are transported across the ocean to other regions.

## 631 **5. Conclusion**

632 We developed a new hybrid algorithm to support the rapid classification of  
633 aerosol types by building an aerosol optical database for global AERONET sites. This  
634 hybrid algorithm is a complex aerosol-type processing algorithm that effectively  
635 integrates machine learning and density clustering algorithms. Additionally, this  
636 algorithm is not limited by the amount of data and improves the accuracy of aerosol-  
637 type classification. On investigating the aerosol types at specific sites with dominant  
638 aerosols, we observed that different sites contained one or more aerosol types, with  
639 the composition of some specific dominant aerosol sites being more complex than that  
640 of others. The new algorithm showed a higher accuracy than that shown by algorithms  
641 used in previous studies in identifying aerosol types at specific sites, particularly in  
642 distinguishing between U/I and mixed-fine aerosols. Finally, the recognition results of  
643 the new hybrid algorithm were closer to the baseline CRI, confirming that the new

644 hybrid algorithm is better than the density-clustering algorithm. On investigating the  
645 aerosol types at global sites across the continents using the new algorithm, we  
646 observed the dominance of different types of aerosols at different sites, and the  
647 composition of these could be logically and effectively attributed to the geographical  
648 location, energy consumption structure, meteorological conditions and activities  
649 happening at the respective sites.

650 In this study, the existing aerosol type identification algorithm was improved  
651 using global ground-based AERONET optical property parameter data, and the spatial  
652 distribution characteristics of global aerosol types were analyzed, which impacted  
653 aerosol radiation research and optical thickness inversion accuracy. Additionally, the  
654 presumption of spherical dust aerosols in the Mie scattering model diverges from their  
655 actual non-spherical nature in the environment, introducing potential inaccuracies.  
656 The optical database's precision, therefore, necessitates further refinement. Future  
657 advancements could involve adopting more potent machine learning techniques, such  
658 as advanced algorithms beyond the current random forest method. Meanwhile, multi-  
659 source satellite data and reanalysis products can be incorporated into aerosol-type  
660 identification. Ultimately, this study will provide support for the identification and  
661 control of air pollution sources.

#### 662 **Author contributions**

663 **Feng Zhang** designed the study. **Xiaoli Wei** analyzed the results and wrote the  
664 original draft. **Qian Cui** engaged in data processing, manuscript editing, and  
665 restructuring. **Leiming Ma** revised the paper and gave constructive suggestions.  
666 **Wenwen Li** gave constructive comments on the paper. **Peng Liu** revised the paper.  
667 All authors contributed to the study.

#### 668 **Competing interests**

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



670 **Acknowledgments**

671 This work was supported by the National Key R&D Program  
672 (2021YFB3900401), the National Natural Science Foundation of China (42075125  
673 and 42105081) and Science and Technology Foundation of Shanghai (23ZR1454100)

674 **References**

- 675 Van G. R., der W., Randerson, J. T., Giglio, L., Collatz, G. J., Mu, M., Kasibhatla, P. S., Morton, D. C.,  
676 Defries, R. S., Jin, Y., and Van Leeuwen, T. T.: Global fire emissions and the contribution of  
677 deforestation, savanna, forest, agricultural, and peat fires (1997–2009), *Atmos. Chem. Phys.*, 10,  
678 11707–11735, <https://doi.org/10.5194/acp-10-11707-2010>, 2010.
- 679 Bahadur, R., Praveen, P. S., Xu, Y., and Ramanathan, V.: Solar absorption by elemental and brown  
680 carbon determined from spectral observations, *Proc. Natl. Acad. Sci. U. S. A.*, 109, 17366–17371,  
681 <https://doi.org/10.1073/pnas.1205910109>, 2012.
- 682 Bian, Yuxuan et al.: Development and Validation of a CCD-Laser Aerosol Detective System for  
683 Measuring the Ambient Aerosol Phase Function., *Atmospheric measurement techniques*, 10  
684 (6),2313–2322. <https://doi.org/10.5194/amt-10-2313>, 2017
- 685 Boselli, A., Caggiano, R., Cornacchia, C., Madonna, F., Mona, L., Macchiato, M., Pappalardo, G., and  
686 Trippetta, S.: Multi year sun-photometer measurements for aerosol characterization in a Central  
687 Mediterranean site, *Atmos. Res.*, 104–105, 98–110, <https://doi.org/10.1016/j.atmosres.2011.08.002>,  
688 2012.
- 689 Breiman: Random forests, *Machine Learning*, 45(1), 5–32, <https://doi.org/10.1023/A:1010933404324>,  
690 2001.
- 691 Che, H., Bing, Q., Zhao, H., Xia, X., and Zhang, X.: Aerosol optical properties and direct radiative  
692 forcing based on measurements from the China Aerosol Remote Sensing Network (CARSNET) in  
693 eastern China, *Atmos. Chem. Phys.*, 18, 405–425, <https://doi.org/10.5194/acp-18-405-2018>, 2018.
- 694 Choi, W., Lee, H., and Park, J.: A first approach to aerosol classification using space-borne  
695 measurement data: Machine learning-based algorithm and evaluation, *Remote Sens.*, 13, 1–21,  
696 <https://doi.org/10.3390/rs13040609>, 2021a.
- 697 Choi, W., Lee, H., Kim, D., and Kim, S.: Improving spatial coverage of satellite aerosol classification  
698 using a random forest model, *Remote Sens.*, 13 (7):1268. <https://doi.org/10.3390/rs13071268>,2021b.
- 699 Dubovik, O. and King, M. D.: A flexible inversion algorithm for retrieval of aerosol optical properties  
700 from Sun and sky radiance measurements, *J. Geophys. Res. Atmos.*, 105, 20673–20696,  
701 <https://doi.org/10.1029/2000JD900282>, 2000.
- 702 Dubovik, O., Holben, B., Eck, T. F., Smirnov, A., Kaufman, Y. J., King, M. D., Tanré, D., and Slutsker,  
703 I.: Variability of absorption and optical properties of key aerosol types observed in worldwide  
704 locations, *J. Atmos. Sci.*, 59, 590–608, <https://doi.org/10.1175/1520-0469>, 2002.
- 705 Eck, T. F., Holben, B. N., Reid, J. S., Dubovik, O., Smirnov, A., O’Neill, N. T., Slutsker, I., and Kinne,  
706 S.: Wavelength dependence of the optical depth of biomass burning, urban, and desert dust aerosols,  
707 *J. Geophys. Res. Atmos.*, 104, 31333–31349, <https://doi.org/10.1029/1999JD900923>, 1999.
- 708 Elham Ghasemifar.:Climatology of aerosol types and their vertical distribution over Iran using  
709 CALIOP dataset during 2007–2021,*Remote Sensing Applications: Society and Environment*,32,  
710 101053, 2352-9385,<https://doi.org/10.1016/j.rsase.2023.101053>.2023.

711 Fernandez-Delgado, M., Cernadas, E., Barro, S., and Amorim, D.: Do we Need Hundreds of Classifiers  
712 to Solve Real World Classification Problems?, *J. Mach. Learn. Res.*, 15, 3133–3181,  
713 <https://dl.acm.org/doi/10.5555/2627435.2697065>, 2014.

714 Fu, Q., Thorsen, T.J., Su, J., Ge, J., & Huang, J.: Test of Mie-based single-scattering properties of non-  
715 spherical dust aerosols in radiative flux calculations. *Journal of Quantitative Spectroscopy &  
716 Radiative Transfer*, 110, 1640-1653. <https://doi.org/10.1016/j.jqsrt.2009.03.010>, 2009

717 Giles, D. M., Holben, B. N., Eck, T. F., Sinyuk, A., Smirnov, A., Slutsker, I., Dickerson, R. R.,  
718 Thompson, A. M., and Schafer, J. S.: An analysis of AERONET aerosol absorption properties and  
719 classifications representative of aerosol source regions, *J. Geophys. Res. Atmos.*, 117, 1–16,  
720 <https://doi.org/10.1029/2012JD018127>, 2012.

721 Hamill, P., Giordano, M., Ward, C., Giles, D., and Holben, B.: An AERONET-based aerosol  
722 classification using the Mahalanobis distance, *Atmos. Environ.*, 140, 213–233,  
723 <https://doi.org/10.1016/j.atmosenv.2016.06.002>, 2016.

724 Kalapureddy, M. C. R., Kaskaoutis, D. G., Ernest Raj, P., Devara, P. C. S., Kambezidis, H. D.,  
725 Kosmopoulos, P. G., and Nastos, P. T.: Identification of aerosol type over the Arabian Sea in the  
726 premonsoon season during the Integrated Campaign for Aerosols, Gases and Radiation Budget  
727 (ICARB), *J. Geophys. Res. Atmos.*, 114, 1–12, <https://doi.org/10.1029/2009JD011826>, 2009.

728 Kaskaoutis, D. G., Kharol, S. K., Sinha, P. R., Singh, R. P., Badarinath, K., Mehdi, W., and Sharma, M.:  
729 Contrasting aerosol trends over South Asia during the last decade based on MODIS observations,  
730 *Atmos. Meas. Tech. Discuss.*, 4, 5275–5323, <https://doi.org/10.5194/amtd-4-5275-2011>, 2011.

731 Kiehl, J. T. and Briegleb, B. P.: The relative roles of sulfate aerosols and greenhouse gases in climate  
732 forcing, *Science (80-. )*, 260, 311–314, <http://dx.doi.org/10.1126/science.260.5106.311>, 1993.

733 Kumar, K. R., Kang, N., and Yin, Y.: Classification of key aerosol types and their frequency  
734 distributions based on satellite remote sensing data at an industrially polluted city in the Yangtze  
735 River Delta, China, *Int. J. Climatol.*, 38, 320–336, <https://doi.org/10.1002/joc.5178>, 2018.

736 Lee, J., Kim, J., Song, C. H., Kim, S. B., Chun, Y., Sohn, B. J., and Holben, B. N.: Characteristics of  
737 aerosol types from AERONET sunphotometer measurements, *Atmos. Environ.*, 44, 3110–3117,  
738 <https://doi.org/10.1016/j.atmosenv.2010.05.035>, 2010.

739 Levy, R. C., Remer, L. A., Mattoo, S., Vermote, E. F., and Kaufman, Y. J.: Second-generation  
740 operational algorithm: Retrieval of aerosol properties over land from inversion of Moderate  
741 Resolution Imaging Spectroradiometer spectral reflectance, *J. Geophys. Res. Atmos.*, 112,  
742 <https://doi.org/10.1029/2006JD007811>, 2007.

743 Li, K., Bai, K., Ma, M., Guo, J., Li, Z., Wang, G., and Chang, N. Bin: Spatially gap free analysis of  
744 aerosol type grids in China: First retrieval via satellite remote sensing and big data analytics, *ISPRS  
745 J. Photogramm. Remote Sens.*, 193, 45–59, <https://doi.org/10.1016/j.isprsjprs.2022.09.001>, 2022.

746 Lin, J., Zheng, Y., Shen, X., Xing, L., and Che, H.: Global aerosol classification based on aerosol  
747 robotic network (Aeronet) and satellite observation, *Remote Sens.*, 13, 1–23,  
748 <https://doi.org/10.3390/rs13061114>, 2021.

749 Ma Lin.: Measurement of aerosol size distribution function using Mie scattering - Mathematical  
750 considerations., *Journal of aerosol science*, 38(11),1150-1162,  
751 <https://doi.org/10.1016/j.jaerosci.2007.08.003>, 2007.

752 Lopez-Pineiro, A., Cabrera, D., Albarran, A., and Pefia, D.: Influence of two-phase olive mill waste  
753 application to soil on terbuthylazine behaviour and persistence under controlled and field conditions,  
754 *J. Soils Sediments*, 11, 771–782, <https://doi.org/10.1007/s11368-011-0362-3>, 2011.

755 Lu, F., Chen, S., Hu, Z., Han, Z., Alam, K., Luo, H., Bi, H., Chen, J., and Guo, X.: Sensitivity and  
756 uncertainties assessment in radiative forcing due to aerosol optical properties in diverse locations in  
757 China, *Sci. Total Environ.*, 860, 160447, <https://doi.org/10.1016/j.scitotenv.2022.160447>, 2023.

758 Michael, I., Mishchenko, and, Larry, D., and Travis: Light scattering by polydisperse, rotationally  
759 symmetric nonspherical particles: Linear polarization, *J. Quant. Spectrosc. Radiat. Transf.*,  
760 [https://doi.org/10.1016/0022-4073\(94\)90130-9](https://doi.org/10.1016/0022-4073(94)90130-9), 1994.

761 Moraes, C. P. A., Fantinato, D. G., and Neves, A.: Epanechnikov kernel for PDF estimation applied to  
762 equalization and blind source separation, *Signal Processing*, 189, 108251,  
763 <https://doi.org/10.1016/j.sigpro.2021.108251>, 2021.

764 Nandan, R., Ratnam, M.V., Kiran, V.R., Madhavan, B.L., & Naik, D.N.: Estimation of Aerosol  
765 Complex Refractive Index over a tropical atmosphere using a synergy of in-situ measurements.  
766 *Atmospheric Research*, 257, 105625, <https://doi.org/10.1016/J.ATMOSRES.2021.105625>,  
767 2021 Nicolae, D., Vasilescu, J., Talianu, C., Biniotoglou, I., Nicolae, V., Andrei, S., and Antonescu, B.:  
768 A neural network aerosol-typing algorithm based on lidar data, *Atmos. Chem. Phys.*, 18, 14511–  
769 14537, <https://doi.org/10.5194/acp-18-14511-2018>, 2018.

770 Omar, A. H., Won, J. G., Winker, D. M., Yoon, S. C., Dubovik, O., and McCormick, M. P.:  
771 Development of global aerosol models using cluster analysis of Aerosol Robotic Network  
772 (AERONET) measurements, *J. Geophys. Res. D Atmos.*, 110, 1–14,  
773 <https://doi.org/10.1029/2004JD004874>, 2005.

774 Pace, G., di Sarra, A., Meloni, D., Piacentino, S., and Chamard, P.: Aerosol optical properties at  
775 Lampedusa (Central Mediterranean). 1. Influence of transport and identification of different aerosol  
776 types, *Atmos. Chem. Phys.*, 6, 697–713, <https://doi.org/10.5194/acp-6-697-2006>, 2006.

777 Papadakis, G. Z., Megaritis, A. G., and Pandis, S. N.: Effects of olive tree branches burning emissions  
778 on PM<sub>2.5</sub> concentrations, *Atmos. Environ.*, 112, 148–158,  
779 <https://doi.org/10.1016/j.atmosenv.2015.04.014>, 2015.

780 Pathak, B., Bhuyan, P. K., Gogoi, M., and Bhuyan, K.: Seasonal heterogeneity in aerosol types over  
781 Dibrugarh-North-Eastern India, *Atmos. Environ.*, 47, 307–315,  
782 <https://doi.org/10.1016/j.atmosenv.2011.10.061>, 2012.

783 Pawar, G. V., Devara, P. C. S., and Aher, G. R.: Identification of aerosol types over an urban site based  
784 on air-mass trajectory classification, *Atmos. Res.*, 164–165, 142–155,  
785 <https://doi.org/10.1016/j.atmosres.2015.04.022>, 2015.

786 Puxbaum, H., Caseiro, A., Sánchez-Ochoa, A., Kasper-Giebl, A., Claeys, M., Gelencsér, A., Legrand,  
787 M., Preunkert, S., and Pio, C.: Levoglucosan levels at background sites in Europe for assessing the  
788 impact of biomass combustion on the European aerosol background, *J. Geophys. Res.*, 112, D23S05,  
789 <https://doi.org/10.1029/2006JD008114>, 2007.

790 Quirantes, Arturo et al.: Extinction-related Angström exponent characterization of submicrometric  
791 volume fraction in atmospheric aerosol particles., *Atmospheric Research*, 228(D24), 270–280,  
792 <https://doi.org/10.1016/j.atmosres.2019.06.009>, 2019

793 Ramanathan, V., Crutzen, P. J., Lelieveld, J., Mitra, A. P., Althausen, D., Anderson, J., Andreae, M. O.,  
794 Cantrell, W., Cass, G. R., and Chung, C. E.: Indian Ocean Experiment: An integrated analysis of the  
795 climate forcing and effects of the great Indo-Asian haze, *J. Geophys. Res. Atmos.*, 106,  
796 <https://doi.org/10.1029/2001JD900133>, 2001.

797 Raut, J. C. and Chazette, P.: Radiative budget in the presence of multi-layered aerosol structures in the  
798 framework of AMMA SOP-0, *Atmos. Chem. Phys.*, 8, 6839–6864, [35](https://doi.org/10.5194/acp-8-</a></p>
</div>
<div data-bbox=)

799 6839-2008, 2008.

800 Reddy LA, Glover TA, Dudek CM, Alperin A, Wiggs NB, Bronstein B.: A randomized trial examining  
801 the effects of paraprofessional behavior support coaching for elementary students with disruptive  
802 behavior disorders: Paraprofessional and student outcomes. *J Sch Psychol.* 2022 Jun;92:227-245.  
803 <https://doi.org/10.1016/j.jsp.2022.04.002>, 2022.Redemann, J., Turco, R. P., Liou, K. N., Russell, P.  
804 B., Bergstrom, R. W., Schmid, B., Hobbs, P. V, Hartley, W. S., Ismail, S., and Ferrare, R. A.:  
805 Retrieving the vertical structure of the effective aerosol complex index of refraction from a  
806 combination of aerosol in situ and remote sensing measurements during TARFOX, *J. Geophys. Res.*,  
807 105( D8), 9949– 9970, doi:10.1029/1999JD901044,2000.

808 Remer, L. A., Tanré, D., and Kaufman, Y. J.: Algorithm for remote sensing of tropospheric aerosol from  
809 MODIS: Collection 005, 2009.

810 Rosenblatt, M.: Remarks on Some Nonparametric Estimates of a Density Function, Remarks on Some  
811 Nonparametric Estimates of a Density Function. In: Davis, R., Lii, KS., Politis, D. (eds) Selected  
812 Works of Murray Rosenblatt. Selected Works in Probability and Statistics. Springer, New York, NY.  
813 [https://doi.org/10.1007/978-1-4419-8339-8\\_13](https://doi.org/10.1007/978-1-4419-8339-8_13), 2011.

814 Sheridan, P. J., Delene, D. J., and Ogren, J. A.: Four Years of Continuous Surface Aerosol  
815 Measurements from the DOE / ARM Southern Great Plains CART Site, 1–8,  
816 <https://doi.org/10.1029/2001JD000785>, 2001.

817 Shin, S. K., Tesche, M., Noh, Y., and Müller, D.: Aerosol-type classification based on AERONET  
818 version 3 inversion products, *Atmos. Meas. Tech.*, 12, 3789–3803, [https://doi.org/10.5194/amt-12-](https://doi.org/10.5194/amt-12-3789-2019)  
819 3789-2019, 2019.

820 Siomos, N., Fountoulakis, I., Natsis, A., Drosoglou, T., and Bais, A.: Automated aerosol classification  
821 from spectral UV measurements using machine learning clustering, *Remote Sens.*, 12, 1–18,  
822 <https://doi.org/10.3390/rs12060965>, 2020.

823 Tanré, D., Kaufman, Y. J., Holben, B. N., Chatenet, B., Karnieli, A., Lavenu, F., Blarel, L., Dubovik, O.,  
824 Remer, L. A., and Smirnov, A.: Climatology of dust aerosol size distribution and optical properties  
825 derived from remotely sensed data in the solar spectrum, *J. Geophys. Res. Atmos.*, 106, 18205–  
826 18217, <https://doi.org/10.1029/2000JD900663>, 2001.

827 Tong, H., Lakey, P. S. J., Arangio, A. M., Socorro, J., Kampf, C. J., Berkemeier, T., Brune, W. H.,  
828 Pöschl, U., and Shiraiwa, M.: Reactive oxygen species formed in aqueous mixtures of secondary  
829 organic aerosols and mineral dust influencing cloud chemistry and public health in the Anthropocene,  
830 *Faraday Discuss.*, 200, 251–270, <https://doi.org/10.1039/c7fd00023e>, 2017.

831 Wang J, Liu Y, Chen L, Liu Y, Mi K, Gao S, Mao J, Zhang H, Sun Y, Ma Z.: Validation and calibration  
832 of aerosol optical depth and classification of aerosol types based on multi-source data over China.  
833 *Sci Total Environ.* 2023 Dec 10;903:166603. doi: 10.1016/j.scitotenv.2023.

834 Wu, Y., Li, J., Xia, Y., Deng, Z., Tao, J., Tian, P., Gao, Z., Xia, X., and Zhang, R.: Size-resolved  
835 refractive index of scattering aerosols in urban Beijing: A seasonal comparison, *Aerosol Sci.*  
836 *Technol.*, 55, 1070–1083, <https://doi.org/10.1080/02786826.2021.1924357>, 2021.

837 Yang, M., Howell, S. G., Zhuang, J., and Huebert, B. J.: Attribution of aerosol light absorption to black  
838 carbon, brown carbon, and dust in China - Interpretations of atmospheric measurements during  
839 EAST-AIRE, *Atmos. Chem. Phys.*, 9, 2035–2050, <https://doi.org/10.5194/acp-9-2035-2009>, 2009.

840 Yokelson, R. J., Urbanski, S. P., Atlas, E. L., Toohey, D. W., Alvarado, E. C., Crounse, J. D., Wennberg,  
841 P. O., Fisher, M. E., Wold, C. E., and Campos, T. L.: Emissions from forest fires near Mexico City ,  
842 *Atmos. Chem. Phys.*, 7, 5569–5584, <https://doi.org/10.5194/acp-7-5569-2007>, 2007.

843 Yousefi, R., Wang, F., Ge, Q., and Shaheen, A.: Long-term aerosol optical depth trend over Iran and  
844 identification of dominant aerosol types, *Sci. Total Environ.*, 722,  
845 <https://doi.org/10.1016/j.scitotenv.2020.137906>, 2020.

846 Zhang, L. and Li, J.: Variability of major aerosol types in China classified using AERONET  
847 measurements, *Remote Sens.*, 11, <https://doi.org/10.3390/rs11202334>, 2019.

848 Zhao, G., Li, F., & Zhao, C.: Determination of the refractive index of ambient aerosols. *Atmospheric*  
849 *Environment*, 240, 117800. <https://doi.org/10.1016/j.atmosenv.2020.117800>,2020

850

851