

1 Title: Evaluation of a Hierarchical Agglomerative Clustering Method Applied to WIBS
2 Laboratory Data for Improved Discrimination of Biological Particles by Comparing Data
3 Preparation Techniques

4
5 NICOLE SAVAGE^{1#}, J Alex Huffman¹

6 ¹ *University of Denver, Department of Chemistry and Biochemistry, Denver, USA*

7 [#] *Now at Aerosol Devices, Inc.*

8
9 *Correspondence to: J. Alex Huffman (alex.huffman@du.edu)*

10
11 Running Title: Evaluation of clustering applied to WIBS bioaerosol data

12
13 Keywords: Clustering, Thresholding, Ward's linkage, Bioaerosols, Fluorescence, Laboratory
14 characterization

15
16
17 **Abstract**

18 Hierarchical agglomerative clustering (HAC) analysis has been successfully applied to
19 several sets of ambient data (e.g. Crawford et al., 2015; Robinson et al., 2013) and with respect
20 to standardized particles in the laboratory environment (Ruske et al., 2017; Ruske et al., 2018).
21 Here we show for the first time a systematic application of HAC to a comprehensive set of
22 laboratory data collected for many individual particle types using the Wideband Integrated
23 Bioaerosol Sensor (WIBS-4A) (Savage et al., 2017). The impact of ratio of particle
24 concentrations on HAC results was investigated, showing that clustering quality can vary
25 dramatically as a function of ratio. Six strategies for particle pre-processing were also compared,
26 concluding that using raw fluorescence intensity (without normalizing to particle size) and
27 inputting all data in logarithmic bins consistently produced the highest quality results for the
28 particle types analyzed. A total of 23 one-on-one matchups of individual particles types were
29 investigated. Results showed cluster misclassification of <15% for 12 of 17 numerical
30 experiments using one biological and one non-biological particle type each. Inputting
31 fluorescence data using a baseline + 3 σ threshold produced lower misclassification than when
32 inputting either all particles (without fluorescence threshold) or a baseline + 9 σ threshold. Lastly,
33 six numerical simulations of mixtures of four to seven components were analyzed using HAC.
34 These results show that a range of 12-24% of fungal clusters were consistently misclassified by
35 inclusion of a mixture of non-biological materials, whereas bacteria and diesel soot were each
36 able to be separated with nearly 100% efficiency. The study gives significant support to the
37 application of clustering analysis to data from commercial UV-LIF instruments being commonly
38 used for bioaerosol research across the globe and provides practical tools that will improve
39 clustering results within scientific studies as a part of diverse research disciplines.

40 1. Introduction

41 Particles of biological origin, or bioaerosols, make up a substantial fraction of atmospheric
42 aerosol and have the potential to influence environmental processes and to negatively impact
43 human health (Després et al., 2012; Douwes et al., 2003; Fröhlich-Nowoisky et al., 2016;
44 Shiraiwa et al., 2017). In order to understand the impact bioaerosols, such as pollen, spores, and
45 bacteria, play on various systems, it is important to be able to identify and characterize these
46 biological particles in the atmosphere. One common method for the detection of bioaerosols is
47 ultraviolet laser/light-induced fluorescence (UV-LIF), because it can provide particle detection in
48 near real-time and at high particle size resolution (Fennelly et al., 2017; Huffman and Santarpia,
49 2017; Sodeau and O'Connor, 2016). Many commercial UV-LIF instruments have become
50 available for bioaerosol detection, but all of these techniques are challenged with the need to
51 differentiate between small differences in fluorescence properties in order to identify and
52 quantify biological aerosols from non-biological material. Recently commercialized instruments
53 show improved ability to discriminate between particle types, for example by utilizing multiple
54 excitation sources or other particle data (e.g. size and shape). UV-LIF techniques are inherently
55 limited, however, by the broad nature of fluorescence spectra and so instruments face a
56 ubiquitous problem of poor selectivity between particle types. By applying improved data
57 thresholding and particle classification techniques, particle characterization can be further
58 improved, but important limitations still remain (Hernandez et al., 2016; Huffman et al., 2012;
59 Perring et al., 2015; Savage et al., 2017; Toprak and Schnaiter, 2013; Wright et al., 2014). One
60 strategy to improving quality of differentiation between particles types has been to collect full,
61 resolved emission spectra, each at multiple excitation wavelengths. This can lead to high
62 instrumental purchase cost, and such instruments have not been widely applied or
63 commercialized (Huffman et al., 2016; Kiselev et al., 2013; Pan et al., 2009b; Ruske et al., 2017;
64 Swanson and Huffman, 2018). Most commercial UV-LIF instruments for bioaerosol detection
65 utilize 1-2 excitation wavelengths and integrate fluorescence signals into a small number of
66 emission bands. To extend the improvements in particle classification for these commercial UV-
67 LIF instruments, a number of multivariate analysis techniques have been applied to ambient
68 particle analysis. The most common of these techniques include principal component analysis,
69 factor analysis, and cluster analysis strategies. Classification algorithms, including several
70 clustering techniques in particular, have shown successful results in providing unbiased insights
71 to the classification of bioaerosols (Crawford et al., 2015; Pinnick et al., 2013; Robinson et al.,
72 2013; Swanson and Huffman, 2018).

73 Cluster analysis is a broad class of data mining methods in which data objects placed in the
74 same group (or cluster) are more similar to one another than to those objects placed in other
75 groups. Classification algorithms can be divided into two central models: (1) supervised and (2)
76 unsupervised learning. Both models have associated advantages and disadvantages. Supervised
77 learning methods allow the “training” of data and grouping to better reflect the data observations
78 (Eick et al., 2004; Ruske et al., 2017; Ruske et al., 2018). This type of method enhances (trains)
79 the classification algorithm in that the output groups are pre-determined rather than discovered,
80 as is the case for unsupervised methods. Supervision requires the user to have appropriate
81 starting conditions to put into the model, which are often difficult or impossible to determine.
82 Supervised training methods are also much more time-efficient compared to unsupervised
83 methods, which is important when analyzing ambient datasets where particle counts (individual
84 objects) can be greater than 10^6 (Ruske et al., 2017). In contrast, unsupervised training methods
85 present less bias and can adapt to unique situations, because the resultant clusters are based on

86 models that have not been previously trained. To access some of the advantages of supervised
87 methods, however, it is important to first apply unsupervised models to wide collections of
88 laboratory data of known particle types in order to gain insight on how these models interpret
89 data inputs and to learn how algorithms can best be trained (Ruske et al., 2017).

90 Hierarchical agglomerative clustering (HAC) is an unsupervised learning method that has
91 been most commonly applied for bioaerosol related studies (e.g. Crawford et al., 2016; Crawford
92 et al., 2015; Gosselin et al., 2016; Pan et al., 2009a; Pan et al., 2007; Pinnick et al., 2013; Pinnick
93 et al., 2004; Robinson et al., 2013; Ruske et al., 2017; Ruske et al., 2018). Other unsupervised
94 clustering techniques, such as the k-means clustering method, have shown poor results when
95 applied to ambient data sets because the number of clusters used to represent the data are
96 required a priori, and this information is usually unknown prior to analysis (Ruske et al., 2017).
97 There are several different HAC methods or linkages including: Single, Complete, Average,
98 Weighted, Ward's, Centroid, and Median (Crawford et al., 2015; Müllner, 2013). Ruske et al.
99 (2017) compared a variety of HAC linkages and determined that Ward's linkage had a higher
100 percentage of correctly classifying particles, in comparison to other HAC methods.

101 Recently, Savage et al. (2017) published a comprehensive laboratory study applying the
102 Wideband Integrated Bioaerosol Sensor (WIBS-4A) to a large and diverse set of biological and
103 non-biological aerosol types. Following on that work, the study presented here utilizes those data
104 as inputs to evaluate and challenge the HAC strategy of particle differentiation using the Ward's
105 linkage of unsupervised clustering. Previous HAC studies have focused primarily on (a) the
106 analysis of simple particle standards (i.e. fluorescent microbeads) and (b) clustering of particles
107 from ambient data sets. There have been relatively few published attempts to differentiate
108 between biological particles and interfering particles by clustering methods using controlled
109 laboratory UV-LIF data or to separate different kinds of biological particles from one another.
110 Presented here are results of the HAC method applied to data from a comprehensive WIBS
111 laboratory study showing that clustering can dramatically improve removal of non-biological
112 particle types from data sets if operated under appropriate conditions.

113 114 **2. Experimental and Computational Methods**

115 The WIBS-4A (Droplet Measurement Techniques, Longmont, CO) is a commonly used UV-
116 LIF based instrument for the detection and characterization of biological particles. The
117 instrument collects particles in the size range 0.8 – 20 μm and interrogates them in real-time as
118 particles flow through the path between optical sources. The WIBS collects information about
119 fluorescence intensity in three channels (FL1, FL2, and FL3), particle size, and particle
120 asymmetry for each interrogated particle. The bands of excitation and fluorescence emission are:
121 FL1 ($\lambda_{\text{ex}} = 280 \text{ nm}$, $\lambda_{\text{em}} = 310 - 400 \text{ nm}$), FL2 ($\lambda_{\text{ex}} = 280 \text{ nm}$, $\lambda_{\text{em}} = 420 - 650 \text{ nm}$), and FL3 (λ_{ex}
122 $= 370 \text{ nm}$, $\lambda_{\text{em}} = 420 - 650 \text{ nm}$). The excitation and emission wavelengths chosen for each of the
123 3 fluorescence channels were designed to maximize the information gained about key biological
124 fluorophores present in a broad range of bioparticles (Kaye et al., 2005; Pöhlker et al., 2012).
125 Early generations of UV-LIF bioaerosol spectrometers were often interpreted to be able to detect
126 proteins via channels similar to FL1 and products of active cellular metabolism (i.e. riboflavin
127 and NAD(P)H) via channels similar to FL3, but these approximations are gross simplifications
128 that confound more detailed investigation of particle types. For more information on the design,
129 operation, and calibration of this instrument see e.g. the manuscripts listed here and references
130 therein (Foot et al., 2008; Healy et al., 2012a; Healy et al., 2012b; Hernandez et al., 2016; Kaye
131 et al., 2005; Perring et al., 2015; Robinson et al., 2017; Savage et al., 2017; Stanley et al., 2011).

132 All aerosol materials utilized have been listed previously in Table 2 shown by Savage et al.
133 (2017), where an overview of size and fluorescence properties of particles utilized for this study
134 are also reported. No additional laboratory experiments were performed here beyond the results
135 presented previously.

136 The fluorescence threshold applied to the differentiation of fluorescent from non-fluorescent
137 particles is a key step in UV-LIF data analysis. Traditionally a fluorescence threshold has been
138 determined as the average baseline fluorescence intensity measured in each of the three channels
139 during the forced trigger (FT) mode when no particles are present, plus three times the standard
140 deviation (σ) of that measurement (i.e. $FT + 3\sigma$) (Gabey et al., 2010). Savage et al. (2017) also
141 reported that additional particle discrimination is possible by using $FT + 9\sigma$ as the threshold.
142 Both threshold definitions will be discussed here. After choosing a threshold of minimum
143 fluorescence, the fluorescence characteristics of a particle can be classified into 7 different
144 particle types introduced by Perring et al. (2015) and as summarized in Figure 1 shown by
145 Savage et al. (2017).

146

147 **3. Clustering Strategy**

148 Hierarchical clustering methods work by grouping objects from the bottom up, meaning that
149 each object (particle) starts as its own “cluster,” and clusters are merged together based on
150 similarities until a greatly reduced number of clusters are presented as a final solution. Ward’s
151 method for clustering is among the most popular approaches for HAC and is the only method
152 based on a classical sum-of-squares criterion, minimizing the within-group sum of squares (or
153 variance) (Müllner, 2013). The WIBS-4A used here for data collection provides 5 parameters of
154 information for each individual particle detected (3 fluorescence channels, size and asymmetry
155 factor:AF), resulting in 5 dimensions of data.

156 The clustering analysis was performed using the open-source software R package
157 ‘fastercluster’ (Müllner, 2013) using a Dell Latitude E7450 laptop computer with an Intel®
158 Core™ Processor (i7-5600U CPU @ 2.60 GHz, 16 GB RAM).

159

160 **3.1 Data Preparation**

161 Saturation of fluorescence intensity occurs at 2047 analog-to-digital counts (ADC) for each
162 of the three FL channels in the WIBS-4A, at which point the photomultiplier tube (PMT) reaches
163 its upper limit of detection. A study by Ruske et al. (2017) investigated whether non-fluorescent
164 (in that case, particles below the $FT + 3\sigma$ fluorescence threshold) and/or saturating data points
165 included in the clustering analysis hindered the efficiency of the cluster output. The authors
166 determined that removing both saturating and non-fluorescent particles before HAC analysis
167 resulted in a better clustering performance in terms of correctly classifying ambient particles.
168 The quality of the clustering results is likely to be impacted by types of particles involved and
169 the assumptions placed on those. As shown by Savage et al. (2017), many biological particles
170 present a large fraction that saturate one or more of the fluorescence detectors. Conversely, many
171 non-biological particles present a large fraction of very weakly fluorescent particles with
172 intensity below a given threshold and thus that are classified as non-fluorescent. To limit pre-
173 modification of particle populations before clustering, the only filter applied before clustering
174 was to remove particles smaller than the lower particle size detection limit of the WIBS-4A (0.8
175 μm), similar to Ruske et al. (2017). In contrast, both saturating and non-fluorescent particles
176 were analyzed and the clustering results will be evaluated. Figure 1 outlines the data preparation

177 process, including the conceptual process of normalization, clustering, and validation of data,
178 which is explained in detail below.

179

180 **3.2 Data Normalization**

181 Normalization of the raw data is necessary before executing the clustering algorithm,
182 because data parameters delivered from the instrument are measured on different respective
183 scales. For example, fluorescent intensity values range from 0 to 2047 ADC, size from 0 to ~20
184 μm , and AF from 0 to 100 arbitrary units. Crawford et al. (2015) performed analysis on
185 polystyrene latex spheres (PSLs) using several different normalization techniques, concluding
186 that z-score normalization was the best technique when looking at cluster performance using
187 Ward's linkage for the separation of PSLs. As a result, we utilize the z-score normalization of
188 Ward's linkage HAC for the presented study. By this type of normalization, the mean value of all
189 data points is subtracted from each individual data point, and then each data point is divided by
190 the standard deviation of all points. Standardization using the z-score method compares results to
191 a normal (Gaussian) population, and we have chosen to standardize our variables to a mean of 0
192 and a variance of 1 so that the output variables would be on comparable scales.

193

194 **3.3 HAC Scenarios**

195 Hierarchical agglomerative clustering performs optimally if all variables (1) are independent
196 of one another and (2) can be described well by a normal (Gaussian) distribution (Norusis,
197 2011). To achieve meaningful results from the clustering analysis data values must, therefore, be
198 input into the clustering algorithm with an understanding of how specific preparatory conditions
199 can significantly impact results. To investigate optimal input conditions a total of 6 clustering
200 scenarios were explored, with conditions summarized in Table 1. The impact of two separate
201 variables were explored within these scenarios by varying: (i) whether fluorescence intensity
202 were pre-normalized by particle size and (ii) whether the data values were input after logarithmic
203 transformation to produce a normal distribution.

204 Ambient particle number vs size distributions can often be well approximated by lognormal
205 distributions, although specific groups of particles, including some bacteria, spores, and pollen,
206 may not always exhibit lognormal distribution. Further, fluorescence intensity has been shown to
207 scale with particle size (e.g. Hill et al., 2001; Sivaprakasam et al., 2011). Several previous studies
208 attempted to utilize HAC for ambient lognormally-distributed particle size data (Crawford et al.,
209 2014; Crawford et al., 2015; Robinson et al., 2013), but applied the assumption that particle
210 fluorescence is normally distributed in a group of particles. If this assumption does not hold to be
211 correct, however, weakly fluorescing particles are likely to be grouped into a single cluster based
212 on the high abundance of these particles (Robinson et al., 2013). Scenarios C, D, and E (Table 1)
213 utilize data input to the clustering algorithm after fluorescence intensity was normalized to
214 particle size (by dividing fluorescence intensity value by light scattering signal when a particle
215 interacts with the diode laser beam) in order to explore whether the assumption that laboratory
216 data should be treated like previously explored ambient data sets and not logged. Scenarios B
217 and D take into account the logging of all parameters, producing normal distributions of all
218 variables (AF, particle size, 3 channels of fluorescence). By this process, data values were input
219 into the algorithm as $\log(\text{value})$ without separately binning the points. For comparison, scenarios
220 E and F explore log-spaced distributions of size and AF, while retaining the assumption that the
221 fluorescence output is normally distributed. Scenario A data is neither logged nor normalized.

222 For comparison, Scenario F represents the input conditions that have been used frequently (e.g.
223 Crawford et al., 2015; Ruske et al., 2017).

224

225 **3.4 Cluster Validation**

226 An important feature of HAC is that it provides clusters in an unsupervised manner, and the
227 user must determine the number of clusters that makes physical sense. One useful tool to
228 systematically determine the optimal number of final clusters is the Calinski-Harabasz (CH)
229 index, which uses the interclass-intraclass distance ratio (Liu et al., 2010). For each clustering
230 output the CH index was calculated for cluster solutions with one through ten clusters, and the
231 solution with the highest CH value was generally determined to be the optimal number of
232 clusters. Figure 2 shows an example CH versus cluster number plot for a mixture of *Aspergillus*
233 *niger* fungal spores mixed with diesel soot particles. The curve suggests the optimal result to be a
234 2-cluster solution for this trial, as was generally the case for investigations where two particle
235 types were mixed before clustering. In order to reduce the length and complexity of discussion,
236 analysis of results in Sections 4.1-4.3 was limited to using cluster products only from the 2-
237 cluster solution. In some cases a 3-cluster solution may have produced higher quality results, but
238 these cases were not investigated.

239

240 **4 Results and Discussion**

241 The analysis of clustering quality was performed systematically and with increasing
242 complexity. Section 4.1 utilizes three pairs of particles types to explore the effect of particle ratio
243 and normalization strategies on cluster performance. Using conclusions from this section,
244 Section 4.2 then expands the exploration to 20 additional pairs of particle types. Section 4.3
245 explores the effect of three different fluorescence thresholding strategies on cluster output.
246 Finally, Section 4.4 investigates the ability of HAC analysis to separate particle types from
247 mixed populations of particle types.

248

249 **4.1 Investigating pre-normalization scenarios and particle input ratio**

250 To explore the ability to separate two distinct populations of particles from one another, three
251 different clustering trials are presented in this section as one-on-one match-ups: (1) *Aspergillus*
252 *niger* (fungal spores, F2) vs. NIST diesel soot (S4), (2) *Pseudomonas stutzeri* (bacteria, B3) vs.
253 NIST diesel soot (S4), and (3) *Aspergillus niger* (fungal spores, F2) vs. California sand (mineral
254 dust, D12). These four particle materials were chosen to represent key classes of coarse particles
255 observed in ambient air. For each trial, a subset of particles from each material type was selected
256 randomly for HAC analysis. The clustering process includes: (i) evaluation of cluster
257 performance based on particle assignment and cluster composition, and (ii) visual representations
258 of cluster outputs using particle type classification introduced by Perring et al. (2015). For each
259 of these three trials, the clustering process was run separately using each of the six scenarios A-F
260 described in Table 1. Additionally, while exploring the optimal data pre-processing scenario, the
261 influence that different concentration ratios of particle types could play in the clustering output
262 was also explored. The cluster process for each trial was performed using four different ratios of
263 particles in each particle set including situations with an equal ratio and where the concentration
264 of each particle type was significantly mismatched. In total, this section represents 57 individual
265 clustering experiments (3 trials x 6 scenarios x 3 particle ratios + 3 additional ratio trials)
266 exploring three independent input variables. The results will be utilized to explore many more
267 individual particle type match-ups in the following sections.

268 The first two trials include diesel soot particles, because light-absorbing carbon aerosol are
269 commonly observed in aerosol samples with anthropogenic influence (Bond et al., 2013), and
270 because they can have fluorescence characteristics difficult to distinguish from small biological
271 particles (e.g. Huffman et al., 2010; Pan et al., 2012; Savage et al., 2017; Yu et al., 2016). For
272 example, when excited by photons with a wavelength of 280 nm, diesel soot can be
273 misinterpreted as single bacterial cells using the WIBS, and so we explored here whether the two
274 particle types could be clustered separately (Pöhlker et al., 2012). The three trials include two
275 examples of biological particles, both exhibiting fluorescent properties, but with different
276 excitation-emission characteristics and with different average particle size.

277 The output of the algorithm reports the particle type from which each particle was input in
278 order to evaluate the accuracy of the clustering. The resulting output of each particle with an
279 assigned cluster number is then compared to the originating particle type to determine
280 classification accuracy. Figure 3 summarizes the relative accuracy of individual clustering
281 experiments by representing the percent of particles misclassified with respect to known input
282 identities (blue bar corresponding to correct classification, red bar and overlaid value
283 corresponding to incorrect classification). The clustering process was generally effective for
284 separating particles correctly when two particle types were considered, but results vary widely
285 across the six scenarios. Several previous studies that used HAC to separate particles within an
286 ambient data set assumed that particle fluorescence is already normally distributed (Crawford et
287 al., 2014; Crawford et al., 2015; Robinson et al., 2013). As a result, these previous studies did
288 not normalize fluorescence data and thus used data preparation scenario F in their clustering
289 analysis. For comparison, scenarios B and D were explored to test whether the clustering
290 efficiency would be improved or hindered by fluorescence normalization. Scenarios A and F
291 produced inconsistent results, with some experiments (i.e. 50:50 ratio of fungal spores:diesel)
292 producing misclassification <1.1%, whereas other experiments (i.e. 20:80 ratio of
293 bacterial:diesel) producing misclassification up to 80%. In contrast, scenarios B and D produced
294 consistently more accurate results. Scenario B, in particular, consistently exhibited the most
295 accurate classification of particles for almost every individual experiment. No experiment
296 involving scenario B produced greater than 9% misclassification of particles, regardless of
297 particle input ratio, and most experiments produced results with 0.1 - 3% error. These
298 observations taken together suggest that particle fluorescence properties may not be well
299 described by normal distributions and that normalizing fluorescence data prior to analysis may
300 be more effective.

301 The results of these experiments also highlight how important the ratio of input particles can
302 be. While scenario B was relatively consistent, varying only between 0.1 and 3.8% error for
303 different ratios of the fungal spore versus diesel match-up, other experiments depended strongly
304 on particle ratio. It is clear that the input ratio of particle types cannot be controlled during an
305 ambient study, and so these results suggest that it is important to keep the possibility of varying
306 concentration ratios in mind when interpreting time- or air mass-associated changes in cluster
307 composition or when relaying the relative confidence in clustering results. For the remainder of
308 the discussion, experiments will be limited to a 50:50 ratio following scenario B. In each case the
309 input particles are a random subset taken from the pool of particles in the experimental data. As a
310 result, individual samples selected from the same experiments (i.e. Fig. 4a, Fig 4e) can show
311 slightly different average properties. In some cases (i.e. diesel soot, Fig. 4d) the number of
312 particles originally analyzed was small and so to keep the input particle ratio 50:50 the
313 corresponding particle type was also limited to small numbers.

314 To extend the investigation of particle input ratio, the three match-ups presented in Figure 3
315 were investigated using Scenario B with 1% bioparticles and 99% non-bioparticles in each
316 respective case. In these experiments the bacteria:diesel soot and fungal spores:dust particles
317 separated relatively well (6.6% and 13.5% misclassification, respectively). The fungal
318 spores:diesel soot separation was poor, however, because the diesel soot particles were nearly
319 evenly split into both clusters, and the fungal spore particles were too low in concentration to
320 influence the cluster properties. More investigation is needed to explore how extreme disparities
321 in particle ratio could negatively influence cluster quality in real-world settings.

322 An important tool readily applied to analysis of ambient data is the categorization of particles
323 into 8 fluorescent particle types (Perring et al., 2015). Thus, to further investigate the quality of
324 cluster accuracy, Figure 4 shows inputs and cluster outputs from three clustering experiments
325 stacked as a function of fluorescence particle type and particle size. The top row of Figure 4
326 shows the input data for *Aspergillus niger* and diesel soot (Fig. 4a-b) paired with the outputs of
327 the 2-cluster solution (Fig. 4g-h). It can be seen that both particle materials have predominantly
328 particle type-A characteristics, meaning that they are fluorescent only in channel FL1. The
329 fungal material also presents roughly a third AB (green) and a small minority of non-fluorescent
330 (gray) characteristics. The size distribution of the fungal spores peaks at $\sim 3 \mu\text{m}$, whereas diesel
331 soot peaks at $\sim 1 \mu\text{m}$ in size. While not shown in this plot style, the spores exhibit moderately
332 higher FL1 channel fluorescence, with a median of 543 ADC, whereas diesel soot exhibits a
333 median of 751 ADC in this channel (see Savage et al., 2017; Table 2). Both particle types show
334 almost no fluorescent characteristics in either FL2 or FL3. In summary, the particle distributions
335 are relatively similar in fluorescence particle type and their differences are largely related to
336 particle size, so separation of these particles through Trial 1 was hypothesized to represent a
337 relatively challenging initial exercise. The clustering outputs presented in Figures 4g-h, however,
338 visually highlight the conclusion represented by Figure 3, which is that the particles in this trial
339 separated very well. Cluster 1 was comprised predominantly of fungal particles and presented
340 fluorescence and size traits qualitatively similar to the input fungal particles, whereas cluster 2
341 was comprised predominantly of diesel soot particles. Results from the 50:50 ratio of the
342 scenario B experiments for the other two trials are also shown in the last two rows of Figure 4. In
343 each case, the qualitative properties of the input particles are extremely well represented by the
344 corresponding output cluster, corroborating the conclusion from Figure 3 that the scenario B
345 cases accurately separated the particle groups investigated through these experiments. It is also
346 important to note here that the method of aerosolization for each particle type plays an important
347 role in the observed size distribution and so results involving laboratory particles should be
348 interpreted with this in mind. Observed fluorescence properties, in contrast, are expected to be
349 conserved at a given particle size and intrinsically related to particle composition.

350

351 **4.2 Investigating cluster quality without fluorescence threshold**

352 After concluding that scenario B exhibited the most consistently accurate clustering results
353 using 2-cluster solutions from mixtures comprised of 2 particle type inputs, the analysis was
354 expanded to include a broader range of particle types. Using 50:50 ratios of two types of input
355 particles, prepared using scenario B (leaving fluorescence data un-normalized and forcing all
356 five data parameters into logarithmically spaced bins), 20 new individual experiments were
357 performed. The results of all 23 experiments (3 from Section 4.1 and 20 introduced in Section
358 4.2) are summarized in Table 2 as the percentage of particle misclassification. These trials were
359 chosen to represent a broad range of individual match-ups that might be expected in ambient air.

360 From the original 69 types of particles analyzed by Savage et al. (2017), 14 were used in
361 experiments here: 8 types of non-biological particles and 6 types of biological particles (2 each
362 of fungal spores, bacteria, and pollen species). Supplemental Figure S4 from Savage et al. (2017)
363 shows size distributions stacked by fluorescence particle type for each of the particle species
364 discussed.

365 Table 2a organizes clustering results into three rows, showing misclassification of F2
366 (*Aspergillus niger* fungal spore), B3 (*Pseudomonas stutzeri* bacteria), and P9 (*Phelum pratense*
367 pollen) particles, respectively, with respect to a variety of other particle types represented by
368 table column. Of the 15 cluster experiments between fungal spore or bacteria and non-biological
369 material (top two table rows), only 3 showed misclassification greater than 7.5% (bold text), and
370 7 were less than 3%. The three outliers were: experiment (7) F2 vs BC3 (glyoxal + ammonium
371 sulfate brown carbon aerosol), (8) F2 vs WT (white t-shirt particles), and (14) B3 vs WT.
372 Looking first at experiment (7), F2 particles show A-type fluorescence characteristics and are
373 dominated by a mode between 1.5 and 4 μm . BC3 particles are primarily non-fluorescent <1.5
374 μm , but are primarily A-type between 1.5 and 3 μm , suggesting similar size and fluorescence
375 properties. The white t-shirt particles separated poorly (~41% misclassification) from both the
376 fungal spore and bacterial particles. All three particle types (WT, F2, and B3) exhibit medium
377 fluorescent intensity in the FL1 channel. The poor ability to separate WT from both F2 and B3
378 was surprising, however, given that WT exhibited significantly higher mean fluorescence in each
379 of the FL2 and FL3 channels. As first mentioned by Savage et al. (2017), great care should be
380 taken when interpreting fluorescent particle results from indoor environments where increased
381 concentrations of bleached fibers from clothing, bedding, paper, and cleaning products may be
382 present.

383 While the results show that the spores and bacterial particles investigated could generally be
384 well separated from most potentially interfering non-biological species, the results were much
385 less successful for differentiation from pollen. P9 pollen particles separated poorly in all
386 experiments (versus D12, H2, or P5), with rate of misclassification ranging from 22 to 47%. It is
387 important to keep in mind, however, that the WIBS was operated using a standard gain setting
388 that limits analysis of particle size to below approximately 20 μm . As a result, the WIBS is
389 insensitive to whole pollen grains and so most of the particles observed during pollen
390 experiments are small pollen fragments. Any intact pollen grains that navigate the flow system to
391 be detected are likely to be binned together in the channel representing the largest particles.
392 Clustering results including pollen should be interpreted accordingly. Pollen grains can fragment
393 in ambient air as function of increased relative humidity (Miguel et al., 2006; Suphioglu et al.,
394 1992; Taylor et al., 2004), but the relative ratio of whole/fragmented particles is hard to predict
395 under ambient conditions. Smaller fragments can also exhibit different fluorescent properties
396 than whole grains (Pöhlker et al., 2013). O'Connor et al. (2014) operated a WIBS-4 (Univ.
397 Hertfordshire) at lower gain in order to improve pollen detection efficiency, but these results are
398 not explored directly here.

399 The WIBS instrument is frequently used to differentiate between airborne biological particles
400 and material of non-biological origin. A secondary goal of differentiating more finely between
401 types of biological aerosols is also frequently pursued. To investigate this goal, six additional
402 experiments were conducted by pairing two different types of non-biological particles (Table
403 2b). In contrast to the results shown in Table 2a, the clustering algorithm showed generally poor
404 ability to separate between two biological particle types. Only one of the six experiments
405 resulted in error <15% (F2 vs B3, 10.3% error), whereas error for the other five experiments

406 ranged from 18% to 65%. The worst accuracy was demonstrated by experiments (22) B1 vs B3
407 and experiment (23) P5 vs P9. Both of these experiments attempted to separate between different
408 species of a single particle type (i.e. between two bacteria or two pollen, respectively). Overall,
409 these results suggest that the clustering strategy may be quite useful at aiding the differentiation
410 of biological material from non-biological material, but that separating more finely to quantify
411 differences between types of individual biological particles is significantly more challenging and
412 not likely to be possible in most situations.

413

414 **4.3 Investigating impact of fluorescence thresholding strategy on cluster quality**

415 In previously published studies, removing particles from clustering analysis that exhibited
416 particle fluorescence intensity below the threshold (i.e. non-fluorescent) or at the saturating point
417 improved the efficiency of clustering (Crawford et al., 2015; Ruske et al., 2017). In Sections 4.1-
418 4.2, particles with either of these characteristics were left in the analysis to prevent the
419 underestimation of particles clustered. In this section, however, we investigated whether
420 removing non-fluorescent particles could improve cluster accuracy for the experiments that
421 performed poorly in Section 4.2. Of the 23 trials represented in Table 2, 10 experiments
422 exhibited 15% or greater misclassification and were subjected to further analysis in order to
423 investigate whether using a more discriminating fluorescence thresholding strategy could
424 improve cluster results. In all 10 cases fluorescence saturating particles were retained, and three
425 separate thresholding conditions were compared by: (I) keeping all non-fluorescent and
426 saturating particles, (II) removing non-fluorescent particles by applying a fluorescence threshold
427 of FT baseline + 3σ , and (III) and removing non-fluorescent particles by applying a fluorescence
428 threshold of FT baseline + 9σ . Savage et al. (2017) showed evidence that applying a FT + 9σ
429 improved WBS results by removing a higher fraction of non-biological material from analysis
430 than by applying the more commonly used FT + 3σ , without negatively impacting observations
431 of biological particles. Table 3 shows the percentage of particles misclassified in each of three
432 scenarios investigated here (Table 3a) as well as the number of particles subjected to the
433 clustering algorithm (Table 3b).

434 Each scenario, with exception of the B3 vs B9 experiment (21), shows a decrease in particle
435 misclassification from scenario I (no fluorescence threshold applied) to scenario II (FT + 3σ). In
436 contrast, eight of the ten scenarios *increase* in particle misclassification when raising the
437 fluorescence threshold from 3σ (II) to 9σ (III). The exceptions to this trend are experiments (8)
438 F2 vs WT and (19) F2 vs P9, which show nominal improvement in error (2-4% reduction) with
439 increased threshold. We hypothesize that the 9σ results degrade, in most cases, because the
440 threshold becomes high enough that most weakly fluorescing particles have been removed from
441 analysis. This reduces the ability of the cluster to group into low and high fluorescence
442 categories, and so remaining particles are separated less efficiently. Secondly, removing particles
443 at higher fluorescence thresholds leads to increasingly poor counting statistics, as represented in
444 Table 3b by the number of particles included in each experiment. Overall, these results suggest
445 that inputting particles into the clustering analysis with at least a nominal fluorescence threshold
446 (i.e. FT + 3σ) can improve the clustering results in many cases, however, increasing the
447 threshold further may decrease cluster quality.

448

449 **4.4 Investigating the capability to separate particles in simulations of complex mixtures**

450 To this point, our investigation has focused on a variety of individual match-ups between two
451 distinct particle types. To better simulate real-world scenarios, we computationally simulated six

452 mixtures of particles by pooling existing WIBS data from selected particle types in prescribed
453 ratios. Each simulated mixture was assembled to roughly represent a different hypothetical
454 mixture of particles that might be expected. Also, the particles in each simulated mixture are
455 assumed to be so dilute that any agglomeration is negligible. Table 4 provides an overview of the
456 percentage of each particle type included as well as the total number of particles in the mixture.
457 Mixtures 1 and 2 were simulated arbitrarily to test if a minority (25%) of one type of fungal
458 spores (F2) could be separated from a majority (75%) of a mixture of three different non-
459 biological materials. Mixtures 3 and 4 synthesized arbitrary mixtures of two types of bioaerosol
460 (F2 and B3) with three or five types of non-biological particles, respectively. Mixture 5 was
461 simulated to examine the separation of pollen (P9) from a set of five non-biological particles.
462 Mixture 6 was simulated to be similar to an indoor environment that might have a mixture of
463 biological particles (F2 and B3) with non-biological materials, including bleached fibers (WT).
464 These mixtures are not intended to closely mimic any set of individual ambient conditions, but
465 are rather used as very rough simulations used for discussion and to prompt discussion related to
466 future experiments within the community. In a real-world sampling environment one would also
467 expect a high concentration of non-fluorescent particles as well (e.g. most organic aerosols, sea
468 salt, dusts), but these were generally not sampled as a part of the Savage et al. (2017) study,
469 which focused on fluorescent particles. As a result, relatively non-fluorescent particles like D12
470 and H2 were included here as “fillers” in most mixtures as surrogates for other types of non-
471 fluorescent particles. Clustering analysis was performed using the ratios listed in Table 4, the B
472 scenario of pre-normalization conditions, and filtering non-fluorescent particles below the FT +
473 3σ threshold. In all cases, the number of clusters retrieved after HAC was pre-defined to be the
474 same as the number of particle types input.

475 Cluster results from all six mixtures are summarized in Figure 5. Figure 5 (Part A) shows the
476 number of particles from each type assigned to each cluster, and Parts B and C show results
477 grouped by general particle classification (brown for non-biological and dark green for
478 biological). Overall, the ability of the HAC analysis to separate the biological particles from the
479 non-biological particles was high. In some cases, the quality of separation of one or two
480 biological species from a mixture of non-biological materials was even higher than the 2-
481 material match-ups shown in Sections 4.1-4.3. The two 4-component mixtures showed 22.4%
482 and 14.8% misclassification of fungal spores. In both cases, a small fraction of each of the non-
483 biological materials were mixed into the spore cluster, whereas almost none (1.5% and 0.6%) of
484 the spores were incorrectly mixed into the sum of the non-biological clusters.

485 Mixtures 3 and 4 showed similar misclassification for fungal spores (11.9% and 13.8%,
486 respectively), whereas the bacterial particles clustered with amazing quality. For Mixture 3, no
487 particles other than bacterial particles were grouped into Cluster 1, and only 16 of 213 bacterial
488 particles were assigned to other clusters. For Mixture 4, 135 of 137 particles in Cluster 6 were
489 bacterial in origin and 135 of 142 bacterial particles were assigned to the cluster. The
490 combination of fungal and bacterial particles in Mixtures 3 and 4 resulted in a total of 5.0% and
491 5.3% misclassification of all biological particles.

492 In contrast to the poor separation of pollen from other particle types discussed in Section 4.2,
493 Mixture 5 showed a higher quality of separation between pollen (9.4% misclassified) and the
494 sum of five other non-biological particle types. Lastly, the mixture designed to roughly mimic an
495 indoor environment including white t-shirt particles. In this mixture the WT particles confounded
496 the spore separation, but the bacterial separation was nearly flawless.

497 Another surprising observation from the analysis of these simulated mixtures was that the
498 diesel soot particles (Mixtures 1, 2, 4, and 5) separated into their own cluster in almost all cases
499 with very high quality (1.8%, 2.9%, 0.6%, and 9.4%, respectively, of diesel soot particles
500 misclassified into a different cluster). The quality of separation of bacterial particles and diesel
501 soot (Mixture 4) was especially amazing, given the qualitative similarity of the two particle
502 populations. For example, size-distributions of each particle type show primarily A-type particles
503 with similar mean fluorescent intensity values in FL1, FL2, and FL3 (Savage et al., 2017).
504

505 **5. Conclusions**

506 Application of results from a recent set of systematic laboratory experiments (Savage et al.,
507 2017) by the commonly used hierarchical agglomerative clustering analysis helps to reveal areas
508 where the tool can be used well and other areas where it struggles. First (Section 4.1) it was
509 observed that differing ratios of particle input into the clustering algorithm can produce
510 dramatically different results. It will be important for anyone applying HAC to ambient particle
511 sets where particle ratios are not independently verified to interpret results somewhat loosely. In
512 Section 4.1 the clustering quality of scenario B, where fluorescence intensity was not normalized
513 to particle size and where all input variables were binned into log space, was determined to
514 consistently demonstrate the highest quality results. Further, the ability to the HAC analysis to
515 separate between two groups of individual particle types using no fluorescence threshold
516 (Section 4.2) and comparing three separate threshold strategies (Section 4.3) was shown to be
517 relatively high in many cases, but confounded in others. Lastly, Section 4.4 explored the ability
518 of HAC analysis to separate biological components from more complex mixtures of four to
519 seven types of input particles.

520 A standard fluorescence threshold of $FT + 3\sigma$ has been commonly applied during WIBS
521 analysis to separate between fluorescent and non-fluorescent particles. Savage et al. (2017)
522 concluded that application of a more aggressive threshold strategy ($FT + 9\sigma$) could help
523 discriminate between biological and non-biological particles more successfully in many
524 circumstances, however certain types of interfering, non-biological particle species can still
525 confound WIBS analysis irrespective of the threshold. Here we have investigated an orthogonal
526 strategy to separate particle types by subjecting particles to HAC computer analysis. By
527 comparing the results of the HAC analysis with raw separation based on fluorescence
528 thresholding alone, the HAC analysis can clearly increase quality of differentiation. Interestingly,
529 while Savage et al. (2017) reported that the $FT + 9\sigma$ strategy helped improved differentiation,
530 using the same threshold in conjunction with HAC analysis actually degraded results. We
531 therefore conclude that if HAC analysis is to be performed, the standard $FT + 3\sigma$ threshold is
532 likely to produce the highest quality results, however if HAC is not to be applied that the $FT +$
533 9σ threshold is probably a better choice to enable investigation of biological particles while
534 computationally filtering non-biological particles.

535 The overall message here is that HAC can be applied successfully to differentiate particle
536 types sampled by WIBS instruments and that it is most successful at separating biological
537 species (i.e. fungal spores and bacteria) from non-biological particles. In all cases the HAC
538 method allows separation of particles at least at the order-of-magnitude level, and often with
539 misclassification of $<5\%$. As mentioned by Savage et al. (2017), however, it should always be
540 kept in mind that different instruments may produce slightly different signals due to physical
541 differences between instruments (i.e. fluorescence calibration, tuning, and detector gain
542 sensitivity) and between calibration strategies (Könemann et al., 2018; Robinson et al., 2017).

543 Results here are also generally extendable to other UV-LIF instruments, whether they offer
544 single or many channels of emission spectral resolution, in that the methods of particle pre-
545 preparation and the impact of particle number ratio are likely to relay similar effects on
546 clustering strategy. Subtle differences in particles observed in a real-world environment may also
547 complicate HAC analysis or the extension of results presented here. The UV-LIF community is
548 encouraged to continue laboratory investigations, including detailed interrogation of clustering
549 analytical techniques, to further understand limitations to better differentiating between particles.
550

551 **6. Acknowledgments**

552 The authors acknowledge the University of Denver for financial support from the faculty
553 start-up fund. Nicole Savage acknowledges financial support from the Phillipson Graduate
554 Fellowship at the University of Denver. Martin Gallagher, David Topping, and Simon Ruske in
555 the School of Earth and Environmental Sciences at the University of Manchester are
556 acknowledged for initial discussion regarding clustering strategy. Cathy Durso at the University
557 of Denver Center for Statistics and Visualization is acknowledged for help running clustering
558 algorithms. All contributors to the Savage et al. (2017) paper, in which all experimental data
559 discussed here were originally presented, are acknowledged for their contributions.

560 **7. References**

- 561
- 562 Bond, T. C., Doherty, S. J., Fahey, D. W., Forster, P. M., Berntsen, T., DeAngelo, B. J., Flanner,
563 M. G., Ghan, S., Karcher, B., Koch, D., Kinne, S., Kondo, Y., Quinn, P. K., Sarofim, M. C.,
564 Schultz, M. G., Schulz, M., Venkataraman, C., Zhang, H., Zhang, S., Bellouin, N., Guttikunda,
565 S. K., Hopke, P. K., Jacobson, M. Z., Kaiser, J. W., Klimont, Z., Lohmann, U., Schwarz, J. P.,
566 Shindell, D., Storelvmo, T., Warren, S. G., and Zender, C. S.: Bounding the role of black carbon
567 in the climate system: A scientific assessment, *J. Geophys. Res.-Atmos.*, 118, 5380-5552, 2013.
- 568 Crawford, I., Lloyd, G., Herrmann, E., Hoyle, C. R., Bower, K. N., Connolly, P. J., Flynn, M. J.,
569 Kaye, P. H., Choulaton, T. W., and Gallagher, M. W.: Observations of fluorescent aerosol-cloud
570 interactions in the free troposphere at the High-Altitude Research Station Jungfraujoch,
571 *Atmospheric Chemistry and Physics*, 16, 2273-2284, 2016.
- 572 Crawford, I., Robinson, N. H., Flynn, M. J., Foot, V. E., Gallagher, M. W., Huffman, J. A.,
573 Stanley, W. R., and Kaye, P. H.: Characterisation of bioaerosol emissions from a Colorado pine
574 forest: results from the BEACHON-RoMBAS experiment, *Atmos. Chem. Phys.*, 14, 8559-8578,
575 2014.
- 576 Crawford, I., Ruske, S., Topping, D. O., and Gallagher, M. W.: Evaluation of hierarchical
577 agglomerative cluster analysis methods for discrimination of primary biological aerosol, *Atmos.*
578 *Meas. Tech.*, 8, 4979-4991, 2015.
- 579 Després, V. R., Huffman, J. A., Burrows, S. M., Hoose, C., Safatov, A. S., Buryak, G. A.,
580 Fröhlich-Nowoisky, J., Elbert, W., Andreae, M. O., Pöschl, U., and Jaenicke, R.: Primary
581 Biological Aerosol Particles in the Atmosphere: A Review, *Tellus Series B-Chemical and*
582 *Physical Meteorology*, 64, 15598, 2012.
- 583 Douwes, J., Thorne, P., Pearce, N., and Heederik, D.: Bioaerosol health effects and exposure
584 assessment: Progress and prospects, *Annals of Occupational Hygiene*, 47, 187-200, 2003.
- 585 Eick, C. F., Zeidat, N., and Zhao, Z.: Supervised clustering-algorithms and benefits, 2004, 774-
586 776.
- 587 Fennelly, M. J., Sewell, G., Prentice, M. B., O'Connor, D. J., and Sodeau, J. R.: The Use of
588 Real-Time Fluorescence Instrumentation to Monitor Ambient Primary Biological Aerosol
589 Particles (PBAP), *Atmosphere*, 9, 1, 2017.
- 590 Foot, V. E., Kaye, P. H., Stanley, W. R., Barrington, S. J., Gallagher, M., and Gabey, A.: Low-
591 cost real-time multi-parameter bio-aerosol sensors, *Proceedings of the SPIE - The International*
592 *Society for Optical Engineering*, 7116, 711601, 2008.
- 593 Fröhlich-Nowoisky, J., Kampf, C. J., Weber, B., Huffman, J. A., Pöhlker, C., Andreae, M. O.,
594 Lang-Yona, N., Burrows, S. M., Gunthe, S. S., Elbert, W., Su, H., Hoor, P., Thines, E.,
595 Hoffmann, T., Després, V. R., and Pöschl, U.: Bioaerosols in the Earth system: Climate, health,
596 and ecosystem interactions, *Atmospheric Research*, 182, 346-376, 2016.

597 Gabey, A. M., Gallagher, M. W., Whitehead, J., Dorsey, J. R., Kaye, P. H., and Stanley, W. R.:
598 Measurements and comparison of primary biological aerosol above and below a tropical forest
599 canopy using a dual channel fluorescence spectrometer, *Atmospheric Chemistry and Physics*, 10,
600 4453-4466, 2010.

601 Gosselin, M. I., Rathnayake, C. M., Crawford, I., Pohlker, C., Frohlich-Nowoisky, J., Schmer,
602 B., Despres, V. R., Engling, G., Gallagher, M., Stone, E., Poschl, U., and Huffman, J. A.:
603 Fluorescent bioaerosol particle, molecular tracer, and fungal spore concentrations during dry and
604 rainy periods in a semi-arid forest, *Atmospheric Chemistry and Physics*, 16, 15165-15184, 2016.

605 Healy, D. A., O'Connor, D. J., Burke, A. M., and Sodeau, J. R.: A laboratory assessment of the
606 Waveband Integrated Bioaerosol Sensor (WIBS-4) using individual samples of pollen and fungal
607 spore material, *Atmospheric Environment*, 60, 534-543, 2012a.

608 Healy, D. A., O'Connor, D. J., and Sodeau, J. R.: Measurement of the particle counting
609 efficiency of the "Waveband Integrated Bioaerosol Sensor" model number 4 (WIBS-4), *Journal*
610 *of Aerosol Science*, 47, 94-99, 2012b.

611 Hernandez, M., Perring, A. E., McCabe, K., Kok, G., Granger, G., and Baumgardner, D.:
612 Chamber catalogues of optical and fluorescent signatures distinguish bioaerosol classes,
613 *Atmospheric Measurement Techniques*, 9, 3283-3292, 2016.

614 Hill, S. C., Pinnick, R. G., Niles, S., Fell, N. F., Pan, Y. L., Bottiger, J., Bronk, B. V., Holler, S.,
615 and Chang, R. K.: Fluorescence from airborne microparticles: dependence on size, concentration
616 of fluorophores, and illumination intensity, *Applied Optics*, 40, 3005-3013, 2001.

617 Huffman, D. R., Swanson, B. E., and Huffman, J. A.: A wavelength-dispersive instrument for
618 characterizing fluorescence and scattering spectra of individual aerosol particles on a substrate,
619 *Atmos. Meas. Tech.*, 9, 3987-3998, 2016.

620 Huffman, J. A. and Santarpia, J.: *Online Techniques for Quantification and Characterization of*
621 *Biological Aerosols*. In: *Microbiology of Aerosols*, John Wiley & Sons, Inc., 2017.

622 Huffman, J. A., Sinha, B., Garland, R. M., Snee-Pollmann, A., Gunthe, S. S., Artaxo, P., Martin,
623 S. T., Andreae, M. O., and Poeschl, U.: Size distributions and temporal variations of biological
624 aerosol particles in the Amazon rainforest characterized by microscopy and real-time UV-APS
625 fluorescence techniques during AMAZE-08, *Atmospheric Chemistry and Physics*, 12, 11997-
626 12019, 2012.

627 Huffman, J. A., Treutlein, B., and Pöschl, U.: Fluorescent biological aerosol particle
628 concentrations and size distributions measured with an Ultraviolet Aerodynamic Particle Sizer
629 (UV-APS) in Central Europe, *Atmospheric Chemistry and Physics*, 10, 3215-3233, 2010.

630 Kaye, P. H., Stanley, W. R., Hirst, E., Foot, E. V., Baxter, K. L., and Barrington, S. J.: Single
631 particle multichannel bio-aerosol fluorescence sensor, *Optics Express*, 13, 3583-3593, 2005.

632 Kiselev, D., Bonacina, L., and Wolf, J.-P.: A flash-lamp based device for fluorescence detection
633 and identification of individual pollen grains, *Review of Scientific Instruments*, 84, 2013.

634 Könemann, T., Savage, N. J., Huffman, J. A., and Pöhlker, C.: Characterization of steady-state
635 fluorescence properties of polystyrene latex spheres using off- and online spectroscopic methods,
636 *Atmospheric Measurement Techniques*, 11, 3987-4003, 2018.

637 Liu, Y., Li, Z., Xiong, H., Gao, X., and Wu, J.: Understanding of internal clustering validation
638 measures, 2010, 911-916.

639 Miguel, A. G., Taylor, P. E., House, J., Glovsky, M. M., and Flagan, R. C.: Meteorological
640 influences on respirable fragment release from Chinese elm pollen, *Aerosol Sci. Technol.*, 40,
641 690-696, 2006.

642 Müllner, D.: fastcluster: Fast hierarchical, agglomerative clustering routines for R and Python,
643 *Journal of Statistical Software*, 53, 1-18, 2013.

644 Norusis, M.: Cluster Analysis. In: *IBM SPSS Statistics 19 Guide to Data Analysis*, Norusis &
645 SPSS Inc., 2011.

646 O'Connor, D. J., Healy, D. A., Hellebust, S., Buters, J. T. M., and Sodeau, J. R.: Using the
647 WIBS-4 (Waveband Integrated Bioaerosol Sensor) Technique for the On-Line Detection of
648 Pollen Grains, *Aerosol Sci. Technol.*, 48, 341-349, 2014.

649 Pan, Y.-L., Pinnick, R. G., Hill, S. C., and Chang, R. K.: Particle-Fluorescence Spectrometer for
650 Real-Time Single-Particle Measurements of Atmospheric Organic Carbon and Biological
651 Aerosol, *Environ. Sci. Technol.*, 43, 429-434, 2009a.

652 Pan, Y. L., Huang, H., and Chang, R. K.: Clustered and integrated fluorescence spectra from
653 single atmospheric aerosol particles excited by a 263- and 351-nm laser at New Haven, CT, and
654 Adelphi, MD, *Journal of Quantitative Spectroscopy & Radiative Transfer*, 113, 2213-2221,
655 2012.

656 Pan, Y. L., Pinnick, R. G., Hill, S. C., and Chang, R. K.: Particle-Fluorescence Spectrometer for
657 Real-Time Single-Particle Measurements of Atmospheric Organic Carbon and Biological
658 Aerosol, *Environ. Sci. Technol.*, 43, 429-434, 2009b.

659 Pan, Y. L., Pinnick, R. G., Hill, S. C., Rosen, J. M., and Chang, R. K.: Single-particle laser-
660 induced-fluorescence spectra of biological and other organic-carbon aerosols in the atmosphere:
661 Measurements at New Haven, Connecticut, and Las Cruces, New Mexico, *J. Geophys. Res.-*
662 *Atmos.*, 112, D24S19, 2007.

663 Perring, A. E., Schwarz, J. P., Baumgardner, D., Hernandez, M. T., Spracklen, D. V., Heald, C.
664 L., Gao, R. S., Kok, G., McMeeking, G. R., McQuaid, J. B., and Fahey, D. W.: Airborne
665 observations of regional variation in fluorescent aerosol across the United States, *J. Geophys.*
666 *Res.-Atmos.*, 120, 1153-1170, 2015.

667 Pinnick, R. G., Fernandez, E., Rosen, J. M., Hill, S. C., Wang, Y., and Pan, Y. L.: Fluorescence
668 spectra and elastic scattering characteristics of atmospheric aerosol in Las Cruces, New Mexico,
669 USA: Variability of concentrations and possible constituents and sources of particles in various
670 spectral clusters, *Atmospheric Environment*, 65, 195-204, 2013.

671 Pinnick, R. G., Hill, S. C., Pan, Y. L., and Chang, R. K.: Fluorescence spectra of atmospheric
672 aerosol at Adelphi, Maryland, USA: measurement and classification of single particles
673 containing organic carbon, *Atmospheric Environment*, 38, 1657-1672, 2004.

674 Pöhlker, C., Huffman, J. A., Förster, J.-D., and Pöschl, U.: Autofluorescence of atmospheric
675 bioaerosols: spectral fingerprints and taxonomic trends of pollen, *Atmospheric Measurement
676 Techniques*, 13, 3369-3392, 2013.

677 Pöhlker, C., Huffman, J. A., and Pöschl, U.: Autofluorescence of atmospheric bioaerosols -
678 fluorescent biomolecules and potential interferences, *Atmospheric Measurement Techniques*, 5,
679 37-71, 2012.

680 Robinson, E. S., Gao, R.-S., Schwarz, J. P., Fahey, D. W., and Perring, A. E.: Fluorescence
681 calibration method for single-particle aerosol fluorescence instruments, *Atmospheric
682 Measurement Techniques*, 10, 1755, 2017.

683 Robinson, N. H., Allan, J. D., Huffman, J. A., Kaye, P. H., Foot, V. E., and Gallagher, M.:
684 Cluster analysis of WIBS single-particle bioaerosol data, *Atmospheric Measurement Techniques*,
685 6, 337-347, 2013.

686 Ruske, S., Topping, D. O., Foot, V. E., Kaye, P. H., Stanley, W. R., Crawford, I., Morse, A. P.,
687 and Gallagher, M. W.: Evaluation of machine learning algorithms for classification of primary
688 biological aerosol using a new UV-LIF spectrometer, *Atmospheric Measurement Techniques*,
689 10, 695, 2017.

690 Ruske, S., Topping, D. O., Foot, V. E., Morse, A. P., and Gallagher, M. W.: Machine learning
691 for improved data analysis of biological aerosol using the WIBS, *Atmospheric Measurement
692 Techniques Discussions*, In Review, 2018.

693 Savage, N. J., Krentz, C. E., Könemann, T., Han, T. T., Mainelis, G., Pöhlker, C., and Huffman,
694 J. A.: Systematic characterization and fluorescence threshold strategies for the wideband
695 integrated bioaerosol sensor (WIBS) using size-resolved biological and interfering particles,
696 *Atmos. Meas. Tech.*, 10, 4279-4302, 2017.

697 Shiraiwa, M., Ueda, K., Pozzer, A., Lammel, G., Kampf, C. J., Fushimi, A., Enami, S., Arangio,
698 A. M., Frohlich-Nowoisky, J., Fujitani, Y., Furuyama, A., Lakey, P. S. J., Lelieveld, J., Lucas,
699 K., Morino, Y., Pöschl, U., Takaharna, S., Takami, A., Tong, H. J., Weber, B., Yoshino, A., and
700 Sato, K.: Aerosol Health Effects from Molecular to Global Scales, *Environ. Sci. Technol.*, 51,
701 13545-13567, 2017.

702 Sivaprakasam, V., Lin, H.-B., Huston, A. L., and Eversole, J. D.: Spectral characterization of
703 biological aerosol particles using two-wavelength excited laser-induced fluorescence and elastic
704 scattering measurements, *Optics Express*, 19, 6191-6208, 2011.

705 Sodeau, J. R. and O'Connor, D. J.: Chapter 16 - Bioaerosol Monitoring of the Atmosphere for
706 Occupational and Environmental Purposes. In: *Comprehensive Analytical Chemistry*, de la
707 Guardia, M. and Armenta, S. (Eds.), Elsevier, 2016.

708 Stanley, W. R., Kaye, P. H., Foot, V. E., Barrington, S. J., Gallagher, M., and Gabey, A.:
709 Continuous bioaerosol monitoring in a tropical environment using a UV fluorescence particle
710 spectrometer, *Atmospheric Science Letters*, 12, 195-199, 2011.

711 Suphioglu, C., Singh, M. B., Taylor, P., Knox, R. B., Bellomo, R., Holmes, P., and Puy, R.:
712 Mechanism of grass-pollen-induced asthma, *The Lancet*, 339, 569-572, 1992.

713 Swanson, B. E. and Huffman, J. A.: Development and characterization of an inexpensive single-
714 particle fluorescence spectrometer for bioaerosol monitoring, *Optics Express*, 26, 3646-3660,
715 2018.

716 Taylor, P. E., Flagan, R. C., Miguel, A. G., Valenta, R., and Glovsky, M. M.: Birch pollen
717 rupture and the release of aerosols of respirable allergens, *Clin. Exp. Allergy*, 34, 1591-1596,
718 2004.

719 Toprak, E. and Schnaiter, M.: Fluorescent biological aerosol particles measured with the
720 Waveband Integrated Bioaerosol Sensor WIBS-4: laboratory tests combined with a one year
721 field study, *Atmospheric Chemistry and Physics*, 13, 225-243, 2013.

722 Wright, T. P., Hader, J. D., McMeeking, G. R., and Petters, M. D.: High Relative Humidity as a
723 Trigger for Widespread Release of Ice Nuclei, *Aerosol Sci. Technol.*, 48, i-v, 2014.

724 Yu, X. W., Wang, Z. B., Zhang, M. H., Kuhn, U., Xie, Z. Q., Cheng, Y. F., Poschl, U., and Su,
725 H.: Ambient measurement of fluorescent aerosol particles with a WIBS in the Yangtze River
726 Delta of China: potential impacts of combustion-related aerosol particles, *Atmospheric
727 Chemistry and Physics*, 16, 11337-11348, 2016.
728

729 **Tables**

730

731 Table 1. Six scenarios explored, with varying combinations of pre-analysis treatment. (1)
732 Fluorescence normalization refers to whether fluorescence intensity was normalized to particle
733 size. (2) Variables logged refers to whether data was manipulated to produce a normal
734 distribution.

735

Parameters	A	B	C	D	E	F
1. Fluorescence Normalization	1. No	1. No	1. Yes	1. Yes	1. Yes	1. No
2. Variables Logged	2. No	2. Yes	2. No	2. Yes	2. Yes, only AF/Size variables	2. Yes, only AF/Size variables

736

737 Table 2. Misclassification of 2-cluster solutions for 23 match-ups of two individual particle types
 738 (equal ratio of particle number, B-scenario) computationally combined before clustering
 739 analysis. Misclassification calculated as the sum percentage of particles misclassified in each
 740 cluster divided by the total number of particles. Three biological particle types (F2, B3, P9)
 741 compared separately to (a) non-biological particle materials and (b) biological particle materials.
 742 Particle number input was a subset of total population of particles experimentally analyzed.

(a)

	Non-biological particle materials							
	Diesel soot (Soot 4)	California sand (Dust 2)	Arizona Test Dust (Dust 12)	Suwannee River Humic Acid (HULIS 2)	Methyl- glyoxal + glycine aerosol (Brown carbon 1)	Glyoxal + amm. sulfate aerosol (Brown carbon 3)	White t-shirt (Misc. 2)	Wood smoke (Soot 6)
	S4	D2	D12	H2	BC1	BC3	WT	WS
<i>Aspergillus niger</i> (Fungi 2)	(1) 0.1%	(3) 2.6%	(4) 6.1%	(5) 4.8%	(6) 2.5%	(7) 23.0%	(8) 40.5%	(9) 7.2%
<i>P. stutzeri</i> (Bacteria 3)	(2) 1.2%		(10) 1.9%	(11) 1.2%	(12) 1.3%	(13) 6.1%	(14) 41.7%	(15) 4.7%
<i>Phelum pratense</i> (Pollen 9)			(16) 22.7%	(17) 23.2%				

(b)

	Biological particle materials				
	<i>S. cerevisiae</i> (Fungi 4)	<i>Phelum pratense</i> (Pollen 9)	<i>P. stutzeri</i> (Bacteria 3)	<i>Taxus baccata</i> (Pollen 5)	<i>B. atrophaeus</i> (Bacteria 1)
	F4	P9	B3	P5	B1
<i>Aspergillus niger</i> (Fungi 2)	(18) 27.9%	(19) 36.4%	(20) 10.3%		
<i>P. stutzeri</i> (Bacteria 3)		(21) 18.3%			(22) 65.4%
<i>Phelum pratense</i> (Pollen 9)				(23) 46.8%	

743

744 **Table 3.** Further exploration of 2-cluster solutions for the 10 match-ups of two individual particle
745 types shown in Table 2 with misclassification >15%. Each match-up shown using three separate
746 fluorescence threshold strategies in advance of particle input into cluster algorithm: (I) all
747 particles included (no fluorescence threshold), (II) particles with fluorescence intensity < FT +
748 3 σ removed, and (III) particles with fluorescence intensity < FT + 9 σ removed. (a) Particle
749 misclassification. (b) Total particle number used for clustering experiment.
750

(a)	Percent misclassified	Bio + Non-bio	Input	(7) F2 + BC3	(8) F2 + WT	(14) B3 + WT	(16) P9 + D12	(17) P9 + H2
			(I) All particles	23.0%	40.5%	41.7%	22.7%	23.2%
			(II) Fluor. > FT + 3 σ	10.3%	36.2%	24.3%	19.3%	3.4%
			(III) Fluor. > FT + 9 σ	41.4%	32.6%	31.8%	45.3%	14.0%
	Bio + Bio	Input	(18) F2 + F4	(19) F2 + P9	(21) B3 + P9	(22) B1 + B3	(23) P9 + P5	
		(I) All particles	27.9%	36.4%	18.8%	65.4%	46.8%	
		(II) Fluor. > FT + 3 σ	13.3%	31.0%	20.0%	77.5%	24.9%	
(III) Fluor. > FT + 9 σ		29.0%	28.6%	29.0%	66.7%	33.9%		
(b)	Number of particles	Bio + Non-bio	Input	(7) F2 + BC3	(8) F2 + WT	(14) B3 + WT	(16) P9 + D12	(17) P9 + H2
			(I) All particles	1,959	565	565	10,359	8,902
			(II) Fluor. > FT + 3 σ	1,000	393	393	171	207
			(III) Fluor. > FT + 9 σ	471	319	319	38	37
	Bio + Bio	Input	(18) F2 + F4	(19) F2 + P9	(21) B3 + P9	(22) B1 + B3	(23) P9 + P5	
		(I) All particles	10,000	8,900	10,000	10,000	10,000	
		(II) Fluor. > FT + 3 σ	9,600	8,500	9,800	10,000	10,000	
(III) Fluor. > FT + 9 σ		9,200	8,100	9,700	10,000	7,895		

751

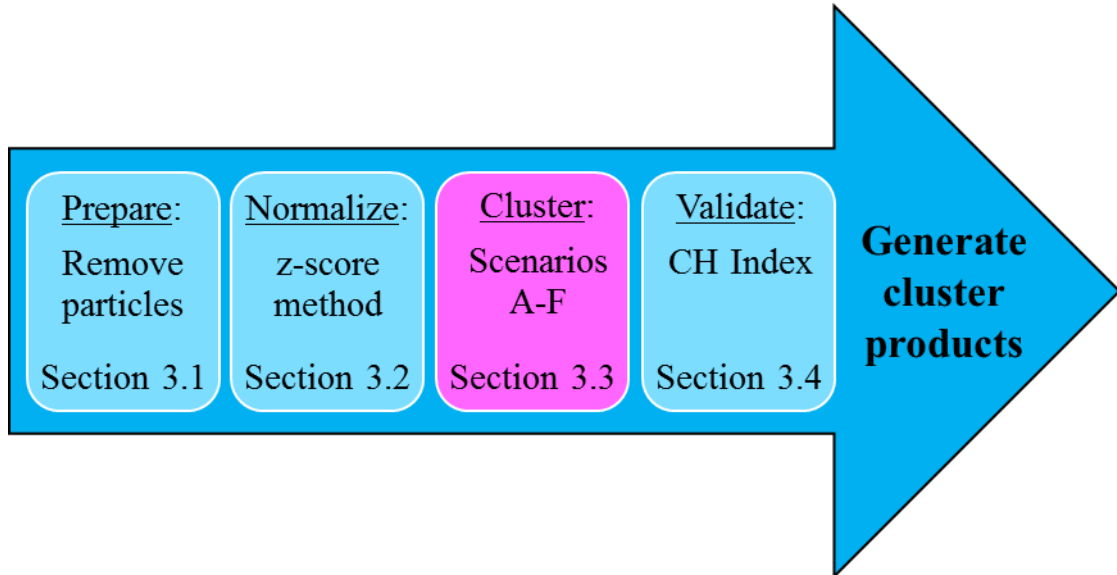
752

753 Table 4. Particle fraction for each type and total particle number used as inputs for simulated
 754 mixtures.
 755

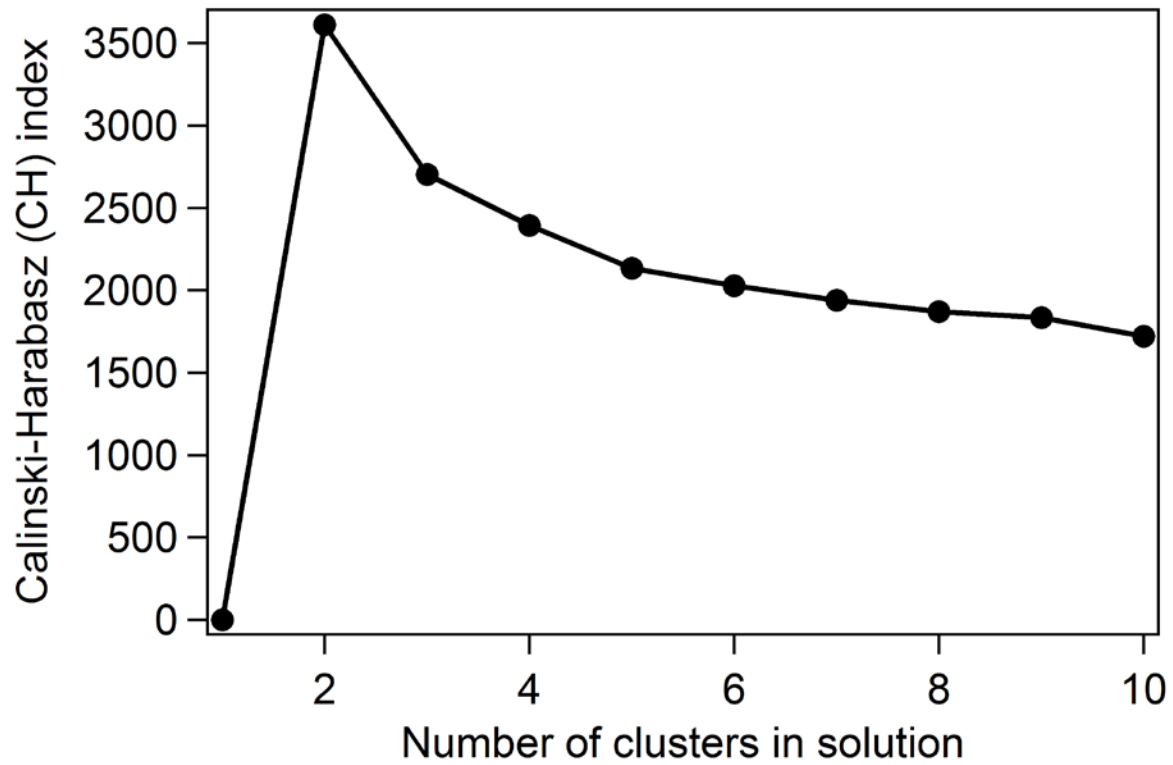
Mixture Number	Mixture Name	F2 <i>Asp. niger</i> (Fungi)	B3 <i>P. stutzeri</i> (Bacteria)	P9 <i>Phelum pretense</i> (Pollen)	S4 Diesel soot	D12 AZ Test Dust	H2 Suwannee River Humic Acid	BC1 Brown Carbon I	WS Wood smoke	WT White t-shirt	Total Particle Number
1	4-Comp. A	25%			25%	25%	25%				680
2	4-Comp. B	25%			25%	25%			25%		680
3	High PBAP	25%	25%			20%	20%	10%			850
4	Low PBAP	12.5%	12.5%		15%	15%	15%	15%	15%		1134
5	Pollen			30%	10%	20%	20%	10%	10%		850
6	Indoor Air	20%	20%			20%	20%			20%	850

756
757

758 **Figures**
759



760
761 Figure 1. Schematic diagram showing the data preparation process resulting in the generated
762 clustering products. Parameters within the pink box are the focus of this manuscript.



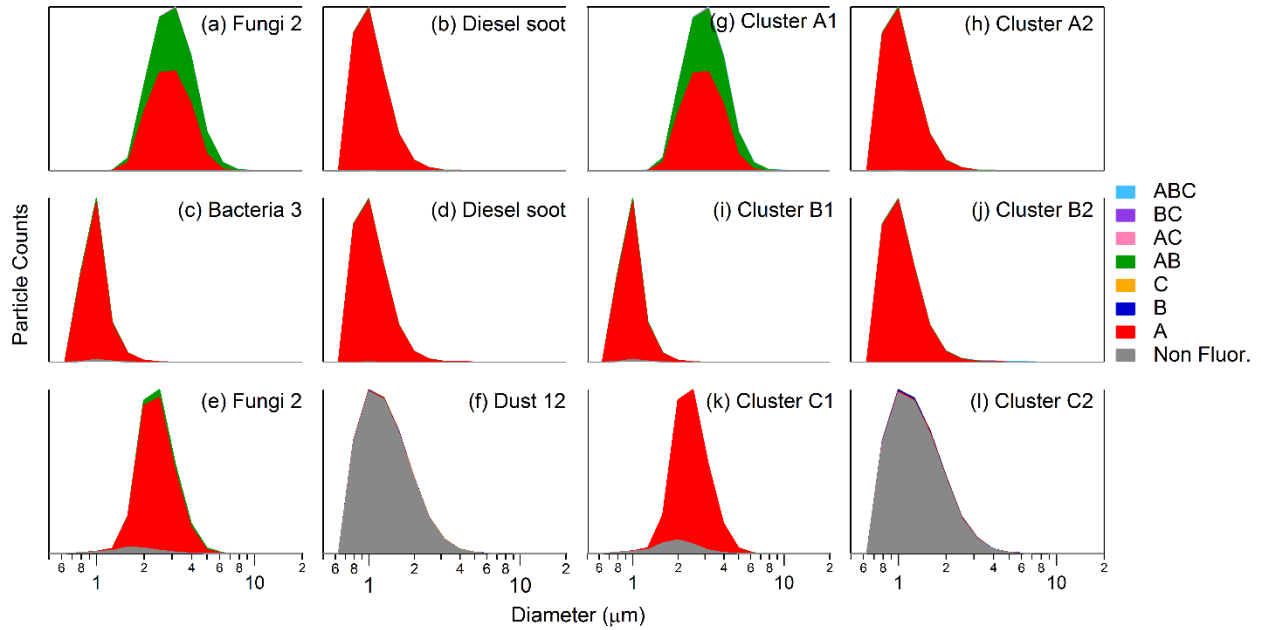
763

764 Figure 2. Example of Calinski-Harabasz Index plot for cluster experiment with input of
765 *Aspergillus niger* and diesel soot (50:50 ratio). Optimal number of clusters is determined by the
766 highest CH value.

	A	B	C	D	E	F
Fungi : Diesel						
50:50 Ratio	1.1	0.9	7.2	4.5	3.6	0.8
80:20 Ratio	64.8	4.1	4.5	2.9	3.8	76.5
20:80 Ratio	2.1	3.8	68.5	6.0	19.5	2.1
Bacteria : Diesel						
50:50 Ratio	50.0	1.2	6.8	4.5	31.6	50.0
80:20 Ratio	0.2	0.2	0.7	1.0	0.9	0.2
20:80 Ratio	80.0	0.3	68.2	0.3	43.7	80.0
Fungi : Dust						
50:50 Ratio	12.7	2.6	24.3	23.5	18.4	30.6
80:20 Ratio	76.6	9.0	20.0	25.4	25.4	29.3
20:80 Ratio	35.9	1.5	55.7	23.4	44.6	58.6

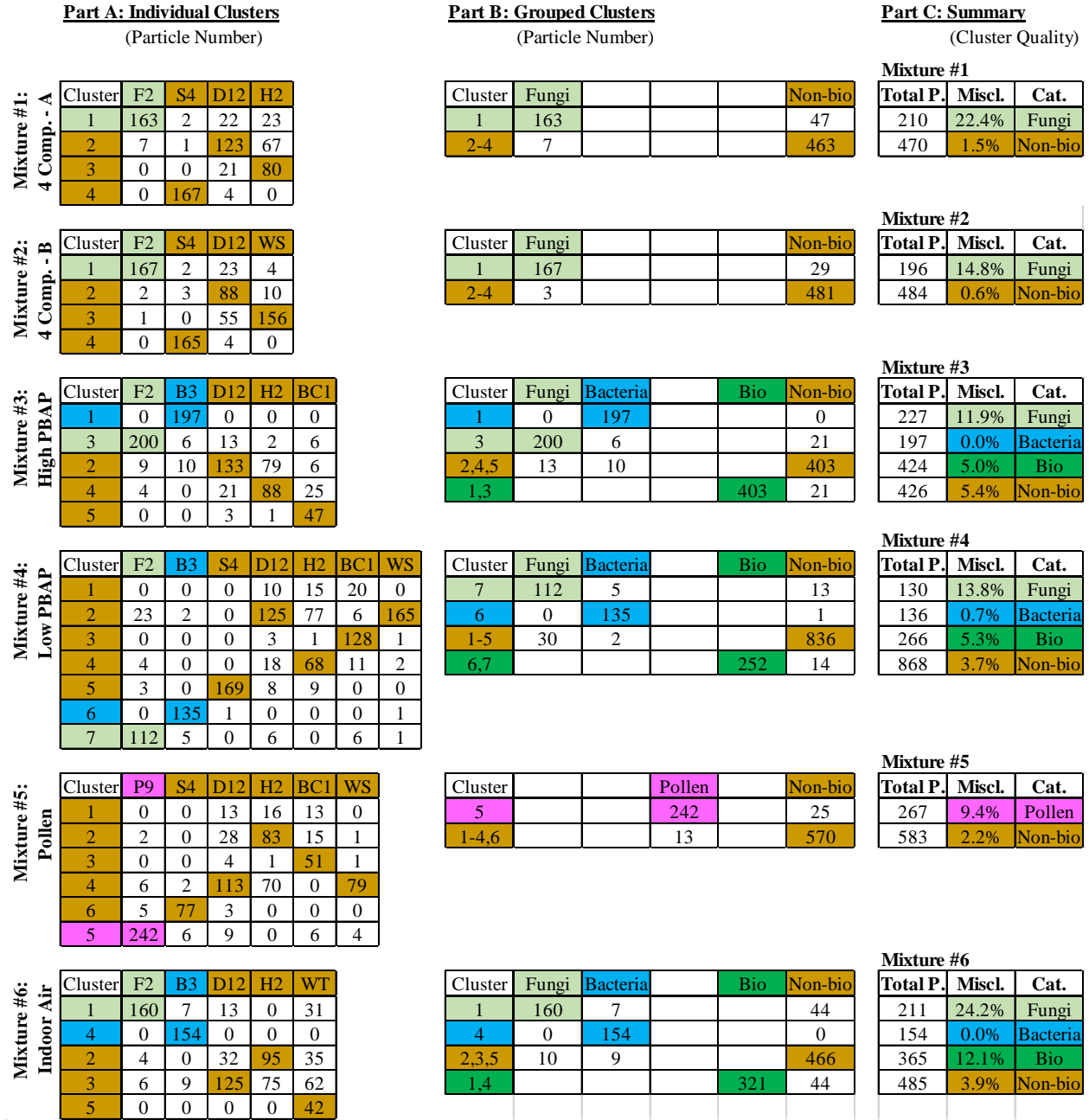
767

768 **Figure 3.** Cluster misclassification shown for three computational combinations of fungal spores
769 (F2), bacteria (B3), diesel soot (S4), and mineral dust (D12). Each combination explored with
770 respect to ratio of input particle number using the scenario B and a 2-cluster solution for each
771 experiment. Scenario letter A-F refers to scenarios summarized in Table 1. Red shaded region
772 (and values) indicates the percent of particles misclassified. Blue shaded region represents the
773 percentage of particles correctly classified.



774
775
776
777
778

Figure 4. Particle type stacked category size distributions for input and output clustering results, using $FT + 3\sigma$ threshold definition. Each experiment (row) shows match-ups of two particle types computationally mixed using 50:50 ratios, scenario B, and 2 cluster solutions. Left two columns show properties of input particles, right two columns show properties of cluster outputs.



779
780
781
782
783
784
785
786
787

Figure 5. Overview of computationally simulated mixtures. Six mixtures shown as groups of rows, with input particle fractions defined in Table 4. Part A (left columns) show particle number retrieved by each individual cluster and categorized by each input particle type. Part B (middle columns) show particle number categorized and grouped by particle classes (i.e. non-biological and biological). Part C (right columns) show misclassification of groups of particles. Colors: light green (fungal spores), blue (bacteria), pink (pollen), dark green (grouped biological), brown (all non-biological).