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 Huffman[1]
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
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). Here we show for
21  the first time a systematic application of HAC to a comprehensive set of laboratory data
22  collected using the wideband integrated bioaerosol sensor (WIBS-4A) (Savage et al., 2017). The
23  impact of particle ratio on HAC results was investigated, showing that clustering quality can
24  vary dramatically as a function of ratio. Six strategies for particle pre-processing were also
25  compared, concluding that using raw fluorescence intensity (without normalizing to particle size)
26  and inputting all data in logarithmic bins consistently produced the highest quality results. A
27  total of 23 one-on-one matchups of individual particles types were investigated. Results showed
28  cluster misclassification of <15% for 12 of 17 analytical experiments using one biological and
29  one non-biological particle type each. Inputting fluorescence data using a baseline + 3σ threshold
30  produced lower misclassification than when inputting either all particles (without fluorescence
31  threshold) or a baseline + 9σ threshold. Lastly, six synthetic mixtures of four to seven
32  components were analyzed. These results show that a range of 12-24% of fungal clusters were
33  consistently misclassified by inclusion of a mixture of non-biological materials, whereas bacteria
34  and diesel soot were each able to be separated with nearly 100% efficiency. The study gives
35  significant support to the application of clustering analysis to data from commercial UV-LIF
36  instruments being commonly used for bioaerosol research across the globe and provides practical
37  tools that will improve clustering results within scientific studies as a part of diverse research
38  disciplines.

## 1. Introduction

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

Cluster analysis is a broad class of data mining methods in which data objects placed in the same group (or cluster) are more similar to one another than to those objects placed in other groups. Clustering techniques can be divided into two central models: (1) supervised and (2) unsupervised learning. Both models have associated advantages and disadvantages. Supervised learning methods allow the "training" of data and grouping to better reflect the data observations (Eick et al., 2004; Ruske et al., 2017). This type of method enhances (trains) the clustering algorithm in that the output cluster classes are pre-determined rather than discovered, as is the case for unsupervised methods. Supervision requires the user to have appropriate starting conditions to put into the model, which are often difficult or impossible to determine. Supervised training methods are also much more time-efficient compared to unsupervised methods, which is important when analyzing ambient datasets where particle counts (individual objects) can be greater than $10^6$ (Ruske et al., 2017). In contrast, unsupervised training methods present less bias and can adapt to unique situations, because the resultant clusters are based on models that have not been previously trained. To access some of the advantages of supervised methods, however,

Atmospheric
Measurement
Techniques
Discussions
Open Access

85  it is critical to first apply unsupervised models to wide collections of laboratory data of known
86  particle types in order to gain insight on how these models interpret data inputs and to learn how
87  algorithms can best be trained (Ruske et al., 2017).
88       Hierarchical agglomerative clustering (HAC) is an unsupervised learning method that has
89  been most commonly applied for bioaerosol related studies (e.g. Crawford et al., 2016; Crawford
90  et al., 2015; Gosselin et al., 2016; Pan et al., 2009a; Pan et al., 2007; Pinnick et al., 2013; Pinnick
91  et al., 2004; Robinson et al., 2013; Ruske et al., 2017). Other unsupervised clustering techniques,
92  such as the k-means clustering method, have shown poor results when applied to ambient data
93  sets because the number of clusters used to represent the data are required a priori, and this
94  information is usually unknown prior to analysis (Ruske et al., 2017). There are several different
95  HAC methods or linkages including: Single, Complete, Average, Weighted, Ward's, Centroid,
96  and Median (Crawford et al., 2015; Müllner, 2013). Ruske et al. (2017) compared a variety of
97  HAC linkages and determined that Ward's linkage had a higher percentage of correctly
98  classifying particles, in comparison to other HAC methods.
99       Recently, Savage et al. (2017) published a comprehensive laboratory study applying the
100 wideband integrated bioaerosol sensor (WIBS-4A) to a large and diverse set of biological and
101 non-biological aerosol types. Following on that work, the study presented here utilizes those data
102 as inputs to evaluate and challenge the HAC strategy of particle differentiation using the Ward's
103 linkage of unsupervised clustering. Previous HAC studies have focused primarily on (a) the
104 analysis of simple particle standards (i.e. fluorescent microbeads) and (b) clustering of particles
105 from ambient data sets. There have been relatively few published attempts to differentiate
106 between biological particles and interfering particles by clustering methods using controlled
107 laboratory UV-LIF data or to separate different kinds of biological particles from one another.
108 Presented here are results of the HAC method applied to data from a comprehensive WIBS
109 laboratory study showing that clustering can dramatically improve removal of non-biological
110 particle types from data sets if operated under appropriate conditions.
111
112     **2.  Experimental and Computing Methods**
113       The WIBS-4A (Droplet Measurement Techniques, Longmont, CO) is a commonly used UV-
114 LIF based instrument for the detection and characterization of biological particles. The
115 instrument collects particles in the size range $0.8 - 20$ μm and interrogates them in real-time as
116 particles flow through the path between optical sources. The WIBS collects 3 channels of
117 fluorescence intensity information (FL1, FL2, and FL3), particle size, and particle asymmetry for
118 each interrogated particle. The excitation and emission wavelengths chosen for each of the 3
119 fluorescence channels were designed to maximize the information gained about key biological
120 fluorophores present in a broad range of bioparticles (Kaye et al., 2005; Pöhlker et al., 2012). For
121 more information on the design, operation, and calibration of this instrument see e.g. the
122 manuscripts listed here and references therein (Foot et al., 2008; Healy et al., 2012a; Healy et al.,
123 2012b; Hernandez et al., 2016; Kaye et al., 2005; Perring et al., 2015; Robinson et al., 2017;
124 Savage et al., 2017; Stanley et al., 2011).
125       All aerosol materials utilized have been listed previously in Table 2 shown by Savage et al.
126 (2017), where an overview of size and fluorescence properties of particles utilized for this study
127 are also reported. No additional laboratory experiments were performed here beyond the results
128 presented previously.
129       The fluorescence threshold applied to the differentiation of fluorescent from non-fluorescent
130 particles is a key step in UV-LIF data analysis. Traditionally a fluorescence threshold has been

Atmospheric
Measurement
Techniques
Discussions

131    determined as the average baseline fluorescence intensity measured in each of three channels
132    during the forced trigger (FT) mode when no particles are present, plus three times the standard
133    deviation ($\sigma$) of that measurement (i.e. FT + 3$\sigma$) (Gabey et al., 2010). Savage et al. (2017) also
134    reported that additional particle discrimination is possible by using FT + 9$\sigma$ as the threshold.
135    Both threshold definitions will be discussed here. After choosing a threshold of minimum
136    fluorescence, the fluorescence characteristics of a particle can be classified into 7 different
137    particle types introduced by Perring et al. (2015) and as summarized in Figure 1 shown by
138    Savage et al. (2017).
139

140    **3.  Clustering Strategy**
141    Hierarchical clustering methods work by grouping objects from the bottom up, meaning that
142    each object (particle) starts as its own "cluster," and clusters are merged together based on
143    similarities until a greatly reduced number of clusters are presented as a final solution. Ward's
144    method for clustering is among the most popular approaches for HAC and is the only method
145    based on a classical sum-of-squares criterion, minimizing the within-group sum of squares (or
146    variance) (Müllner, 2013). The WIBS-4A used here for data collection provides 5 parameters of
147    information for each individual particle detected (3 fluorescence channels, size and asymmetry
148    factor:AF), resulting in 5 dimensions of data.
149    The clustering analysis was performed using the open-source software R package
150    'fastercluster' (Müllner, 2013) using a Dell Latitude E7450 laptop computer with an Intel®
151    Core™ Processor (i7-5600U CPU @ 2.60 GHz, 16 GB RAM).
152

153    **3.1 Data Preparation**
154    Saturation of fluorescence intensity occurs at 2047 analog-to-digital counts (ADC) for each
155    of the three FL channels in the WIBS-4A, at which point the photomultiplier tube (PMT) reaches
156    its upper limit of detection. A study by Ruske et al. (2017) investigated whether non-fluorescent
157    (in that case, particles below the FT + 3$\sigma$ fluorescence threshold) and/or saturating data points
158    included in the clustering analysis hindered the efficiency of the cluster output. The authors
159    determined that removing both saturating and non-fluorescent particles before HAC analysis
160    resulted in a better clustering performance in terms of correctly classifying ambient particles.
161    Their conclusions, however, were based on ambient field data using unknown particles types and
162    did not investigate laboratory-generated particles of known origin. The quality of the clustering
163    results are likely to be impacted by types of particles involved and the assumptions placed on
164    those. As shown by Savage et al. (2017), many biological particles present a large fraction that
165    saturate one or more of the fluorescence detectors. Conversely, many non-biological particles
166    present a large fraction of very weakly fluorescent particles with intensity below a given
167    threshold and thus that are classified as non-fluorescent. To limit pre-modification of particle
168    populations before clustering, the only filter applied before clustering was to remove particles
169    smaller than the lower particle size detection limit of the WIBS-4A (0.8 µm), similar to Ruske et
170    al. (2017). In contrast, both saturating and non-fluorescent particles were retained and the
171    clustering results will be evaluated. Figure 1 outlines the data preparation process, including the
172    conceptual process of normalization, clustering, and validation of data, which will be explained
173    in detail below.
174

Atmospheric
Measurement
Techniques
Discussions

### 3.2 Data Normalization

175
176     Normalization of the raw data is necessary before executing the clustering algorithm,
177 because data parameters delivered from the instrument are measured on different respective
178 scales. For example, fluorescent intensity values range from 0 to 2047 ADC (analog-to-digital
179 counts), size from 0 to ~20 µm, and AF from 0 to 100 arbitrary units. Crawford et al. (2015)
180 performed analysis on polystyrene latex spheres (PSLs) using several different normalization
181 techniques, concluding that z-score normalization is the best technique when looking at cluster
182 performance using Ward's linkage for the separation of PSLs. As a result, we utilize the z-score
183 normalization of Ward's linkage HAC for the presented study. By this type of normalization, the
184 mean value of all data points is subtracted from each individual data point, and then each data
185 point is divided by the standard deviation of all points. Standardization using the z-score method
186 compares results to a normal (Gaussian) population, and it therefore relies on the assumption that
187 input data can be described by a normal distribution (Gordon, 2006).
188
189     **3.3 HAC Scenarios**
190     Hierarchical agglomerative clustering performs optimally if all variables (1) are independent
191 of one another and (2) can be described well by a normal (Gaussian) distribution (Norusis,
192 2011). To achieve meaningful results from the clustering analysis data values must, therefore, be
193 input into the clustering algorithm with a careful understanding of how specific preparatory
194 conditions can significantly impact results. To investigate optimal input conditions a total of 6
195 clustering scenarios were explored, with conditions summarized in Table 1. The impact of two
196 separate variables were explored within these scenarios by varying (i) whether fluorescence
197 intensity were pre-normalized by particle size and (ii) whether the data values were input in
198 logarithmically spaced bins to produce a normal distribution.
199     Ambient particle distributions are well known to exhibit lognormal distributions. Further,
200 fluorescence intensity has been shown to scale with particle size (e.g. Hill et al., 2001;
201 Sivaprakasam et al., 2011). Several previous studies attempted to utilize HAC for ambient
202 lognormally-distributed particle size data (Crawford et al., 2014; Crawford et al., 2015; Robinson
203 et al., 2013), but applied the assumption that particle fluorescence is normally distributed in a
204 group of particles. If this assumption does not hold to be correct, however, weakly fluorescing
205 particles are likely to be grouped into a single cluster based on the high abundance of these
206 particles (Robinson et al., 2013). Scenarios C, D, and E (Table 1) utilize data input to the
207 clustering algorithm after fluorescence intensity was normalized to particle size in order to
208 explore whether the assumption that laboratory data should be treated like previously explored
209 ambient data sets and not logged. Scenarios B and D take into account the logging of all
210 parameters, producing normal distributions of all variables (AF, particle size, 3 channels of
211 fluorescence). For comparison, scenarios E and F explore log-spaced distributions of size and
212 AF, while retaining the assumption that the fluorescence output is normally distributed. Scenario
213 A data is neither logged nor normalized. For comparison, Scenario F represents the input
214 conditions that have been used frequently (e.g. Crawford et al., 2015; Ruske et al., 2017).
215
216     **3.4 Cluster Validation**
217     An important feature of HAC is that it provides clusters in an unsupervised manner, and the
218 user must determine the number of clusters that makes physical sense. One useful tool to
219 systematically determine the optimal number of final clusters is the Calinski-Harabasz (CH)
220 index, which uses the interclass-intraclass distance ratio (Liu et al., 2010). For each clustering

Atmospheric
Measurement
Techniques
Discussions

221  output the CH index was calculated for cluster solutions with one through ten clusters, and the
222  solution with the highest CH value was generally determined to be the optimal number of
223  clusters. Figure 2 shows an example CH versus cluster number plot for a mixture of *Aspergillus*
224  *niger* fungal spores mixed with diesel soot particles. The curve suggests the optimal result to be a
225  2-cluster solution for this trial, as was generally the case for investigations where two particle
226  types were mixed before clustering. In order to reduce the length and complexity of analysis, all
227  cases presented in Sections 4.1-4.3 are products of a 2-cluster solution.

## 4   Results and Discussion

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

### 4.1 Investigating pre-normalization scenarios and particle input ratio

239      To explore the ability to separate two distinct populations of particles from one another, three
240  different clustering trials are presented in this section as one-on-one match-ups: (1) *Aspergillus*
241  *niger* (fungal spores, F2) vs. NIST diesel soot (S4), (2) *Pseudomonas stutzeri* (bacteria, B3) vs.
242  NIST diesel soot (S4), and (3) *Aspergillus niger* (fungal spores, F2) vs. California sand (mineral
243  dust, D12). These four particle materials were chosen to represent key classes of coarse particles
244  observed in ambient air. For each trial, a given number of particles from each material type was
245  placed into a conceptual pool before running through the algorithm to organize clusters. The
246  clustering process includes: (i) evaluation of cluster performance based on particle assignment
247  and cluster composition, and (ii) visual representations of cluster outputs using particle type
248  classification introduced by Perring et al. (2015). For each of these three trials, the clustering
249  process was run separately using each of the six scenarios A-F described in Table 1.
250  Additionally, while exploring the optimal data pre-processing scenario, the influence that
251  different concentration ratios of particle types could play in the clustering output was also
252  explored. The cluster process for each trial was performed using three different ratios of particles
253  in each particle set including an equal ratio (50:50) and situations where the concentration of
254  each particle type was significantly mismatched (80:20 and 20:80). In total, this section
255  represents 54 individual clustering experiments (3 trials x 6 scenarios x 3 particle ratios)
256  exploring three independent input variables. The results will be utilized to explore many more
257  individual particle type match-ups in the following sections.
258      The first two trials include diesel soot particles, because they are commonly observed in
259  almost all atmospheric samples with even minimal anthropogenic influence, and because they
260  have fluorescence characteristics difficult to distinguish from small biological particles (e.g.
261  Huffman et al., 2010; Pan et al., 2012; Savage et al., 2017). For example, when excited by
262  photons with a wavelength of 280 nm, diesel soot can be misinterpreted as single bacterial cells
263  using the WIBS, and so we explored here whether the two particle types could be clustered
264  separately (Pöhlker et al., 2012). The three trials include two examples of biological particles,
265  both exhibiting fluorescent properties, but with different excitation-emission characteristics and
266  with different average particle size.

Atmospheric
Measurement
Techniques
Discussions

267   The output of the algorithm reports the particle type from which each particle was input in
268 order to evaluate the accuracy of the clustering. The resulting output of each particle with an
269 assigned cluster number is then compared to the originating particle type to determine
270 classification accuracy. Figure 3 summarizes the relative accuracy of individual clustering
271 experiments by representing the percent of particles misclassified with respect to known input
272 identities (blue bar corresponding to correct classification, red bar and overlaid value
273 corresponding to incorrect classification). The clustering process was generally effective for
274 separating particles correctly when two particle types were considered, but results vary widely
275 across the six scenarios. Several previous studies that used HAC to separate particles within an
276 ambient data set assumed that particle fluorescence is already normally distributed (Crawford et
277 al., 2014; Crawford et al., 2015; Robinson et al., 2013). As a result, these previous studies did
278 not normalize fluorescence data and thus used data preparation scenario F in their clustering
279 analysis. For comparison, scenarios B and D were explored to test whether the clustering
280 efficiency would be improved or hindered by fluorescence normalization. Scenarios A and F
281 produced inconsistent results, with some experiments (i.e. 50:50 ratio of fungal spores:diesel)
282 producing misclassification <1.1%, whereas other experiments (i.e. 20:80 ratio of
283 bacterial:diesel) producing misclassification >80%. In contrast, scenarios B and D produced
284 consistently more accurate results. Scenario B, in particular, consistently exhibited the most
285 accurate classification of particles for almost every individual experiment. No experiment
286 involving scenario B produced greater than 9% misclassification of particles, regardless of
287 particle input ratio, and most experiments produced results with 0.1 - 3% error. These
288 observations taken together suggest that particle fluorescence properties may not be well
289 described by normal distributions and that normalizing fluorescence data prior to analysis may
290 be more effective.
291   The results of these experiments also highlight how important the ratio of input particles can
292 be. While scenario B was relatively consistent, varying only between 0.1 and 3.8% error for
293 different ratios of the fungal spore versus diesel match-up, other experiments depended strongly
294 on particle ratio. It is clear that the input ratio of particle types cannot be controlled during an
295 ambient study, and so these results suggest that it is important to keep the possibility of varying
296 concentration ratios in mind when interpreting time- or air mass-associated changes in cluster
297 composition or when relaying the relative confidence in clustering results. For the remainder of
298 the discussion, experiments will be limited to a 50:50 ratio following scenario B. In each case the
299 number of input particles represents a random subset taken from the pool of particles in the
300 experimental data. As a result, individual samples selected from the same experiments (i.e. Fig.
301 4a, Fig 4e) can show slightly different average properties. In some cases (i.e. Diesel soot, Fig.
302 4d) the number of particles originally analyzed was small and so to keep the input particle ratio
303 50:50 the corresponding particle type was also limited to small numbers.
304   An important tool readily applied to analysis of ambient data is the categorization of particles
305 into 8 fluorescent particle types (Perring et al., 2015). Thus, to further investigate the quality of
306 cluster accuracy, Figure 4 shows inputs and cluster outputs from three clustering experiments
307 stacked as a function of fluorescence particle type and particle size. The top row of Figure 4
308 shows the input data for *Aspergillus niger* and diesel soot (Fig. 4a-b) paired with the outputs of
309 the 2-cluster solution (Fig. 4g-h). It can be seen that both particle materials have predominantly
310 particle type-A characteristics, meaning that they are fluorescent only in channel FL1. The
311 fungal material also presents roughly a third AB (green) and a small minority of non-fluorescent
312 (gray) characteristics. The size distribution of the fungal spores peaks at ~3 μm, whereas diesel

Atmospheric
Measurement
Techniques
Discussions

Open Access

EGU

313    soot peaks at ~1 μm in size. While not shown in this plot style, the spores exhibit moderately
314    higher FL1 channel fluorescence, with a median of 543 ADC, whereas diesel soot exhibits a
315    median of 751 ADC in this channel (see Savage et al., 2017; Table 2). Both particle types show
316    almost no fluorescent characteristics in either FL2 or FL3. In summary, the particle distributions
317    are relatively similar in fluorescence particle type and their differences are largely related to
318    particle size, so separation of these particles through Trial 1 was hypothesized to represent a
319    relatively challenging initial exercise. The clustering outputs presented in Figures 4g-h, however,
320    visually highlight the conclusion represented by Figure 3, which is that the particles in this trial
321    separated very well. Cluster 1 was comprised predominantly of fungal particles and presented
322    fluorescence and size traits qualitatively similar to the input fungal particles, whereas cluster 2
323    was comprised predominantly of diesel soot particles. Results from the 50:50 ratio of the
324    scenario B experiments for the other two trials are also shown in the last two rows of Figure 4. In
325    each case, the qualitative properties of the input particles are extremely well represented by the
326    corresponding output cluster, corroborating the conclusion from Figure 3 that the scenario B
327    cases accurately separated the particle groups investigated through these experiments.
328
329    **4.2 Investigating cluster quality without fluorescence threshold**
330       After concluding that scenario B exhibited the most consistently accurate clustering results
331    using 2-cluster solutions from mixtures comprised of 2 particle type inputs, the analysis was
332    expanded to include a broader range of particle types. Using 50:50 ratios of two types of input
333    particles, prepared using scenario B (leaving fluorescence data un-normalized and forcing all
334    five data parameters into logarithmically spaced bins), 20 new individual experiments were
335    performed. The results of all 23 experiments (3 from Section 4.1 and 20 introduced in Section
336    4.2) are summarized in Table 2 as the percentage of particle misclassification. These trials were
337    chosen to represent a broad range of individual match-ups that might be expected in ambient air.
338    From the original 69 types of particles analyzed by Savage et al. (2017), 14 were used in
339    experiments here: 8 types of non-biological particles and 6 types of biological particles (2 each
340    of fungal spores, bacteria, and pollen species). Supplemental Figure S4 from Savage et al. (2017)
341    shows size distributions stacked by fluorescence particle type for each of the particle species
342    discussed.
343       Table 2a organizes clustering results into three rows, showing misclassification of F2
344    (*Aspergillus niger* fungal spore), B3 (*Pseudomonas stutzeri* bacteria), and P9 (*Phelum pratense*
345    pollen) particles, respectively, with respect to a variety of other particle types represented by
346    table column. Of the 15 cluster experiments between fungal spore or bacteria and non-biological
347    material (top two table rows), only 3 showed misclassification greater than 7.5% (bold text), and
348    7 were less than 3%. The three outliers were: experiment (7) F2 vs BC3 (glyoxal + ammonium
349    sulfate brown carbon aerosol), (8) F2 vs WT (white t-shirt particles), and (14) B3 vs WT.
350    Looking first at experiment (7), F2 particles show A-type fluorescence characteristics and are
351    dominated by a mode between 1.5 and 4 μm. BC3 particles are primarily non-fluorescent <1.5
352    μm, but are primarily A-type between 1.5 and 3 μm, suggesting similar size and fluorescence
353    properties. The white t-shirt particles separated poorly (~41% misclassification) from both the
354    fungal spore and bacterial particles. All three particle types (WT, F2, and B3) exhibit medium
355    fluorescent intensity in the FL1 channel. The poor ability to separate WT from both F2 and B3
356    was surprising, however, given that WT exhibited significantly higher mean fluorescence in each
357    of the FL2 and FL3 channels. As first mentioned by Savage et al. (2017), great care should be
358    taken when interpreting fluorescent particle results from indoor environments where increased

Atmospheric
Measurement
Techniques
Discussions

359 concentrations of bleached fibers from clothing, bedding, paper, and cleaning products may be
360 present.
361     While the results show that the spores and bacterial particles investigated could generally be
362 well separated from most potentially interfering non-biological species, the results were much
363 less successful for differentiation from pollen. P9 pollen particles separated poorly in all
364 experiments (versus D12, H2, or P5), with rate of misclassification ranging from 22 to 47%. It is
365 important to keep in mind, however, that the WIBS was operated using a standard gain setting
366 that limits analysis of particle size to below approximately 20 μm. As a result, the WIBS is
367 insensitive to whole pollen grains and so most of the particles observed during pollen
368 experiments are small pollen fragments. Any intact pollen grains that navigate the flow system to
369 be detected are likely to be binned together in the channel representing the largest particles.
370 Clustering results including pollen should be interpreted accordingly. Pollen gains can fragment
371 in ambient air as function of increased relative humidity (Miguel et al., 2006; Suphioglu et al.,
372 1992; Taylor et al., 2004), but the relative ratio of whole/fragmented particles is hard to predict
373 under ambient conditions. Smaller fragments can also exhibit different fluorescent properties
374 than whole grains (Pöhlker et al., 2013). O'Connor et al. (2014) operated a WIBS-4 (Univ.
375 Hertfordshire) at lower gain in order to improve pollen detection efficiency, but these results are
376 not explored directly here.
377     The WIBS instrument is frequently used to differentiate between airborne biological particles
378 and material of non-biological origin. A secondary goal of differentiating more finely between
379 types of biological aerosols is also frequently pursued. To investigate this goal, six additional
380 experiments were conducted by pairing two different types of non-biological particles (Table
381 2b). In contrast to the results shown in Table 2a, the clustering algorithm showed generally poor
382 ability to separate between two biological particle types. Only one of the six experiments
383 resulted in error <15% (F2 vs B3, 10.3% error), whereas error for the other five experiments
384 ranged from 18% to 65%. The worst accuracy was demonstrated by experiments (22) B1 vs B3
385 and experiment (23) P5 vs P9. Both of these experiments attempted to separate between different
386 species of a single particle type (i.e. between two bacteria or two pollen, respectively). Overall,
387 these results suggest that the clustering strategy may be quite useful at aiding the differentiation
388 of biological material from non-biological material, but that separating more finely to quantify
389 differences between types of individual biological particles is likely to be significantly more
390 challenging.
391
392 **4.3 Investigating impact of fluorescence thresholding strategy on cluster quality**
393     In previously published studies, removing particles from clustering analysis that exhibited
394 particle fluorescence intensity below the threshold (i.e. non-fluorescent) or at the saturating point
395 improved the efficiency of clustering (Crawford et al., 2015; Ruske et al., 2017). In Sections 4.1-
396 4.2, particles with either of these characteristics were left in the analysis to prevent the
397 underestimation of particles clustered. In this section, however, we investigated whether
398 removing non-fluorescent particles could improve cluster accuracy for the experiments that
399 performed poorly in Section 4.2. Of the 23 trials represented in Table 2, 10 experiments
400 exhibited 15% or greater misclassification and were subjected to further analysis in order to
401 investigate whether using a more discriminating fluorescence thresholding strategy could
402 improve cluster results. In all 10 cases fluorescence saturating particles were retained, and three
403 separate thresholding conditions were compared by: (I) keeping all non-fluorescent and
404 saturating particles, (II) removing non-fluorescent particles by applying a fluorescence threshold

Atmospheric
Measurement
Techniques
Discussions

405     of FT baseline + 3σ, and (III) and removing non-fluorescent particles by applying a fluorescence
406     threshold of FT baseline + 9σ. Table 3 shows the percentage of particles misclassified in each of
407     three scenarios (Table 3a) as well as the number of particles subjected to the clustering algorithm
408     (Table 3b).
409       Each scenario, with exception of the B3 vs B9 experiment (21), shows a decrease in particle
410     misclassification from scenario I (no fluorescence threshold applied) to scenario II (FT + 3σ). In
411     contrast, eight of the ten scenarios *increase* in particle misclassification when raising the
412     fluorescence threshold from 3σ (II) to 9σ (III). The exceptions to this trend are experiments (8)
413     F2 vs WT and (19) F2 vs P9, which show nominal improvement in error (2-4% reduction) with
414     increased threshold. We hypothesize that the 9σ results degrade, in most cases, because the
415     threshold becomes high enough that most weakly fluorescing particles have been removed from
416     analysis. This reduces the ability of the cluster to group into low and high fluorescence
417     categories, and so remaining particles are separated less efficiently. Secondly, removing particles
418     at higher fluorescence thresholds leads to increasingly poor counting statistics, as represented in
419     Table 3b by the number of particles included in each experiment. Overall, these results suggest
420     that inputting particles into the clustering analysis with at least a nominal fluorescence threshold
421     (i.e. FT + 3σ) can improve the clustering results in many cases, however, increasing the
422     threshold further may decrease cluster quality.
423
424 **4.4 Investigating cluster ability to separate complex synthetic mixtures**
425       To this point, our investigation has focused on a variety of individual match-ups between two
426     distinct particle types. To better simulate real-world scenarios, we analytically synthesized six
427     mixtures of particles by pooling existing data from selected particle types in prescribed ratios.
428     Each mixture was synthesized to roughly represent a different hypothetical mixture of particles
429     that might be expected. Table 4 provides an overview of the percentage of each particle type
430     included as well as the total number of particles in the mixture. Mixtures 1 and 2 were
431     synthesized arbitrarily to test if a minority (25%) of one type of fungal spores (F2) could be
432     separated from a majority (75%) of a mixture of three different non-biological materials.
433     Mixtures 3 and 4 synthesized arbitrary mixtures of two types of bioaerosol (F2 and B3) with
434     three or five types of non-biological particles, respectively. Mixture 5 was synthesized to
435     examine the separation of pollen (P9) from a set of five non-biological particles. Mixture 6 was
436     synthesized to simulate an indoor environment that might have a mixture of biological particles
437     (F2 and B3) with non-biological materials, including bleached fibers (WT). These mixtures are
438     not intended to closely mimic any set of individual ambient conditions, but are rather used as
439     very rough synthetic scenarios used for discussion. In a real-world sampling environment one
440     would also expect a high concentration of non-fluorescent particles as well (e.g. most organic
441     aerosols, sea salt, dusts), but these were largely not sampled as a part of the Savage et al. (2017)
442     study, which focused on fluorescent particles. As a result, relatively non-fluorescent particles
443     like D12 and H2 were included here as "fillers" in most mixtures as surrogates for other types of
444     non-fluorescent particles. Clustering analysis was performed using the ratios listed in Table 4,
445     the B scenario of pre-normalization conditions, and filtering non-fluorescent particles below the
446     FT + 3σ threshold. In all cases, the number of clusters retrieved after HAC was the same as the
447     number of particle types input.
448       Cluster results from all six mixtures are summarized in Figure 5. Figure 5 (Part A) shows the
449     number of particles from each type assigned to each cluster, and Parts B and C show results
450     grouped by general particle classification (brown for non-biological and dark green for

Atmospheric
Measurement
Techniques
Discussions

451   biological). Overall, the ability of the HAC analysis to separate the biological particles from the
452   non-biological particles was high. In some cases the quality of separation of one or two
453   biological species from a mixture of non-biological materials was even higher than the 2-
454   material match-ups shown in Sections 4.1-4.3. The two 4-component mixtures showed 22.4%
455   and 14.8% misclassification of fungal spores. In both cases, a small fraction of each of the non-
456   biological materials were mixed into the spore cluster, whereas almost none (1.5% and 0.6%) of
457   the spores were incorrectly mixed into the sum of the non-biological clusters.
458       Mixtures 3 and 4 showed similar misclassification for fungal spores (11.9% and 13.8%,
459   respectively), whereas the bacterial particles clustered with amazing quality. For Mixture 3, no
460   particles other than bacterial particles were grouped into Cluster 1, and only 16 of 213 bacterial
461   particles were assigned to other clusters. For Mixture 4, 135 of 137 particles in Cluster 6 were
462   bacterial in origin and 135 of 142 bacterial particles were assigned to the cluster. The
463   combination of fungal and bacterial particles in Mixtures 3 and 4 resulted in a total of 5.0% and
464   5.3% misclassification of all biological particles.
465       In contrast to the poor separation of pollen from other particle types discussed in Section 4.2,
466   Mixture 5 showed a higher quality of separation between pollen (9.4% misclassified) and the
467   sum of five other non-biological particle types. Lastly, the mixture designed to roughly mimic an
468   indoor environment including white t-shirt particles. In this mixture the WT particles confounded
469   the spore separation, but the bacterial separation was nearly flawless.
470       Another surprising observation from the analysis of these synthetic mixtures was that the
471   diesel soot particles (Mixtures 1, 2, 4, and 5) separated into their own cluster in almost all cases
472   with very high quality (1.8%, 2.9%, 0.6%, and 9.4%, respectively, of diesel soot particles
473   misclassified into a different cluster). The quality of separation of bacterial particles and diesel
474   soot (Mixture 4) was especially amazing, given the qualitative similarity of the two particle
475   populations. For example, size-distributions of each particle type show primarily A-type particles
476   with similar mean fluorescent intensity values in FL1, FL2, and FL3 (Savage et al., 2017).

477
478   **5. Conclusions**
479       Application of results from a recent set of systematic laboratory experiments (Savage et al.,
480   2017) by the commonly used hierarchical agglomerative clustering analysis helps to reveal areas
481   where the tool can be used well and other areas where it struggles. First (Section 4.1) it was
482   observed that differing ratios of particle input into the clustering algorithm can produce
483   dramatically different results. It will be important for anyone applying HAC to ambient particle
484   sets where particle ratios are not independently verified to interpret results somewhat loosely. In
485   Section 4.1 the clustering quality of scenario B, where fluorescence intensity was not normalized
486   to particle size and where all input variables were binned into log space, was determined to
487   consistently demonstrate the highest quality results. Further, the ability to the HAC analysis to
488   separate between two groups of individual particle types using no fluorescence threshold
489   (Section 4.2) and comparing three separate threshold strategies (Section 4.3) was shown to be
490   relatively high in many cases, but confounded in others. Lastly, Section 4.4 explored the ability
491   of HAC analysis to separate biological components from more complex mixtures of four to
492   seven types of input particles.
493       A standard fluorescence threshold of FT + 3σ has been commonly applied during WIBS
494   analysis to separate between fluorescent and non-fluorescent particles. Savage et al. (2017)
495   concluded that application of a more aggressive threshold strategy (FT + 9σ) could help
496   discriminate between biological and non-biological particles more successfully in many

Atmospheric
Measurement
Techniques
Discussions

Open Access

EGU

497    circumstances, however certain types of interfering, non-biological particle species can still
498    confound WIBS analysis irrespective of the threshold. Here we have investigated an orthogonal
499    strategy to separate particle types by subjecting particles to HAC computer analysis. By
500    comparing the results of the HAC analysis with raw separation based on fluorescence
501    thresholding alone, the HAC analysis can clearly increase quality of differentiation. Interestingly,
502    while Savage et al. (2017) reported that the FT + 9σ strategy helped improved differentiation,
503    using the same threshold in conjunction with HAC analysis actually degraded results. We
504    therefore conclude that if HAC analysis is to be performed, the standard FT + 3σ threshold is
505    likely to produce the highest quality results, however if HAC is not to be applied that the FT +
506    9σ threshold is the most likely to reduce a large fraction of non-biological particles.
507        The overall message here is that HAC can be applied successfully to differentiate particle
508    types sampled by WIBS instruments and that it is most successful at separating biological
509    species (i.e. fungal spores and bacteria) from non-biological particles. In all cases the HAC
510    method allows separation of particles at least at the order-of-magnitude level, and often with
511    misclassification of <5%. As mentioned by Savage et al. (2017), however, it should always been
512    kept in mind that different instruments may produce slightly different signals due to physical
513    differences (i.e. fluorescence calibration, tuning, and detector gain sensitivity) and that results
514    here are only generally extendable to other UV-LIF instruments(Robinson et al., 2017). Subtle
515    differences in particles observed in a real-world environment may complicate HAC analysis or
516    the extension of results presented here. The UV-LIF community is encouraged to continue
517    laboratory investigations, including detailed interrogation of clustering analytical techniques, to
518    further understand limitations to better differentiating between particles.
519
520    **6. Acknowledgments**

Atmospheric
Measurement
Techniques
Discussions

# 7. References

Crawford, I., Lloyd, G., Herrmann, E., Hoyle, C. R., Bower, K. N., Connolly, P. J., Flynn, M. J., Kaye, P. H., Choularton, T. W., and Gallagher, M. W.: Observations of fluorescent aerosol-cloud interactions in the free troposphere at the High-Altitude Research Station Jungfraujoch, Atmospheric Chemistry and Physics, 16, 2273-2284, 2016.

Crawford, I., Robinson, N. H., Flynn, M. J., Foot, V. E., Gallagher, M. W., Huffman, J. A., Stanley, W. R., and Kaye, P. H.: Characterisation of bioaerosol emissions from a Colorado pine forest: results from the BEACHON-RoMBAS experiment, Atmos. Chem. Phys., 14, 8559-8578, 2014.

Crawford, I., Ruske, S., Topping, D. O., and Gallagher, M. W.: Evaluation of hierarchical agglomerative cluster analysis methods for discrimination of primary biological aerosol, Atmos. Meas. Tech., 8, 4979-4991, 2015.

Després, V. R., Huffman, J. A., Burrows, S. M., Hoose, C., Safatov, A. S., Buryak, G. A., Fröhlich-Nowoisky, J., Elbert, W., Andreae, M. O., Pöschl, U., and Jaenicke, R.: Primary Biological Aerosol Particles in the Atmosphere: A Review, Tellus Series B-Chemical and Physical Meteorology, 64, 15598, 2012.

Douwes, J., Thorne, P., Pearce, N., and Heederik, D.: Bioaerosol health effects and exposure assessment: Progress and prospects, Annals of Occupational Hygiene, 47, 187-200, 2003.

Eick, C. F., Zeidat, N., and Zhao, Z.: Supervised clustering-algorithms and benefits, 2004, 774-776.

Fennelly, M. J., Sewell, G., Prentice, M. B., O'Connor, D. J., and Sodeau, J. R.: The Use of Real-Time Fluorescence Instrumentation to Monitor Ambient Primary Biological Aerosol Particles (PBAP), Atmosphere, 9, 1, 2017.

Foot, V. E., Kaye, P. H., Stanley, W. R., Barrington, S. J., Gallagher, M., and Gabey, A.: Low-cost real-time multi-parameter bio-aerosol sensors, Proceedings of the SPIE - The International Society for Optical Engineering, 7116, 711601, 2008.

Fröhlich-Nowoisky, J., Kampf, C. J., Weber, B., Huffman, J. A., Pöhlker, C., Andreae, M. O., Lang-Yona, N., Burrows, S. M., Gunthe, S. S., Elbert, W., Su, H., Hoor, P., Thines, E., Hoffmann, T., Després, V. R., and Pöschl, U.: Bioaerosols in the Earth system: Climate, health, and ecosystem interactions, Atmospheric Research, 182, 346-376, 2016.

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

Gordon, S.: The Normal Distribution. University of Syndey, 2006.

Atmospheric
Measurement
Techniques
Discussions

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

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

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

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

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

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

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

586    Huffman, J. A., Sinha, B., Garland, R. M., Snee-Pollmann, A., Gunthe, S. S., Artaxo, P., Martin,
587    S. T., Andreae, M. O., and Poeschl, U.: Size distributions and temporal variations of biological
588    aerosol particles in the Amazon rainforest characterized by microscopy and real-time UV-APS
589    fluorescence techniques during AMAZE-08, Atmospheric Chemistry and Physics, 12, 11997-
590    12019, 2012.

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

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

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

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

Atmospheric
Measurement
Techniques
Discussions

600   Miguel, A. G., Taylor, P. E., House, J., Glovsky, M. M., and Flagan, R. C.: Meteorological
601   influences on respirable fragment release from Chinese elm pollen, Aerosol Sci. Technol., 40,
602   690-696, 2006.

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

605   Norusis, M.: Cluster Analysis. In: IBM SPSS Statistics 19 Guide to Data Analysis, Norusis &
606   SPSS Inc., 2011.

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

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

613   Pan, Y. L., Huang, H., and Chang, R. K.: Clustered and integrated fluorescence spectra from
614   single atmospheric aerosol particles excited by a 263-and 351-nm laser at New Haven, CT, and
615   Adelphi, MD, Journal of Quantitative Spectroscopy & Radiative Transfer, 113, 2213-2221,
616   2012.

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

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

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

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

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

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

638    Pöhlker, C., Huffman, J. A., and Pöschl, U.: Autofluorescence of atmospheric bioaerosols -
639    fluorescent biomolecules and potential interferences, Atmospheric Measurement Techniques, 5,
640    37-71, 2012.

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

644    Robinson, N. H., Allan, J. D., Huffman, J. A., Kaye, P. H., Foot, V. E., and Gallagher, M.:
645    Cluster analysis of WIBS single-particle bioaerosol data, Atmospheric Measurement Techniques,
646    6, 337-347, 2013.

647    Ruske, S., Topping, D. O., Foot, V. E., Kaye, P. H., Stanley, W. R., Crawford, I., Morse, A. P.,
648    and Gallagher, M. W.: Evaluation of machine learning algorithms for classification of primary
649    biological aerosol using a new UV-LIF spectrometer, Atmospheric Measurement Techniques,
650    10, 695, 2017.

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

655    Shiraiwa, M., Ueda, K., Pozzer, A., Lammel, G., Kampf, C. J., Fushimi, A., Enami, S., Arangio,
656    A. M., Frohlich-Nowoisky, J., Fujitani, Y., Furuyama, A., Lakey, P. S. J., Lelieveld, J., Lucas,
657    K., Morino, Y., Poschl, U., Takaharna, S., Takami, A., Tong, H. J., Weber, B., Yoshino, A., and
658    Sato, K.: Aerosol Health Effects from Molecular to Global Scales, Environ. Sci. Technol., 51,
659    13545-13567, 2017.

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

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

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

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

671    Swanson, B. E. and Huffman, J. A.: Development and characterization of an inexpensive single-
672    particle fluorescence spectrometer for bioaerosol monitoring, Optics Express, 26, 3646-3660,
673    2018.

674    Taylor, P. E., Flagan, R. C., Miguel, A. G., Valenta, R., and Glovsky, M. M.: Birch pollen
675    rupture and the release of aerosols of respirable allergens, Clin. Exp. Allergy, 34, 1591-1596,
676    2004.

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

680    Wright, T. P., Hader, J. D., McMeeking, G. R., and Petters, M. D.: High Relative Humidity as a
681    Trigger for Widespread Release of Ice Nuclei, Aerosol Sci. Technol., 48, i-v, 2014.
682

Atmospheric
Measurement
Techniques
Open Access

Discussions

EGU

683 **Tables**

684

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

689

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

690

691   Table 2. Misclassification of 2-cluster solutions for 23 match-ups of two individual particle types
692   (equal ratio of particle number, B-scenario). Misclassification calculated as the sum percentage
693   of particles misclassified in each cluster divided by the total number of particles. Three
694   biological particle types (F2, B3, P9) compared separately to (a) non-biological particle materials
695   and (b) biological particle materials. Particle number input was a subset of total population of
696   particles experimentally analyzed.

(a)

| | Non-biological particle materials | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | Diesel soot (Soot 4) S4 | California sand (Dust 2) D2 | Arizona Test Dust (Dust 12) D12 | Suwannee River Humic Acid (HULIS 2) H2 | Methyl-glyoxal + glycine aerosol (Brown carbon 1) BC1 | Glyoxal + amm. sulfate aerosol (Brown carbon 3) BC3 | White t-shirt (Misc. 2) WT | Wood smoke (Soot 6) WS |
| *Aspergillus niger* (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% |
| *P. stutzeri* (Bacteria 3) | *(2)* *1.2%* | | (10) 1.9% | (11) 1.2% | (12) 1.3% | (13) 6.1% | (14) **41.7%** | (15) 4.7% |
| *Phelum pretense* (Pollen 9) | | | (16) **22.7%** | (17) **23.2%** | | | | |

(b)

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

697

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

(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 + $3\sigma$ | 10.3% | 36.2% | 24.3% | 19.3% | 3.4% |
| | (III) Fluor. > FT + $9\sigma$ | 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 + $3\sigma$ | 13.3% | 31.0% | 20.0% | 77.5% | 24.9% |
| | (III) Fluor. > FT + $9\sigma$ | 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 + $3\sigma$ | 1,000 | 393 | 393 | 171 | 207 |
| | (III) Fluor. > FT + $9\sigma$ | 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 + $3\sigma$ | 9,600 | 8,500 | 9,800 | 10,000 | 10,000 |
| | (III) Fluor. > FT + $9\sigma$ | 9,200 | 8,100 | 9,700 | 10,000 | 7,895 |

Atmospheric
Measurement
Techniques
Discussions
Open Access
EGU

707 Table 4. Particle fraction for each type and total particle number used as inputs for synthetic
708 mixtures.
709

| Mixture Number | Mixture Name | F2 Asp. niger (Fungi) | B3 P. stutzeri (Bacteria) | P9 Phelum pretense (Pollen) | S4 Diesel soot | D12 AZ Test Dust | H2 Suwannee River Humic Acid | BC1 Brown Carbon 1 | 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 |

710
711

Atmospheric
Measurement
Techniques
Open Access
Discussions
EGU

712    **Figures**
713



714
715    Figure 1. Schematic diagram showing the data preparation process resulting in the generated
716    clustering products. Parameters within the pink box are the focus of this manuscript.

Atmospheric
Measurement
Techniques
Discussions

Open Access



717

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

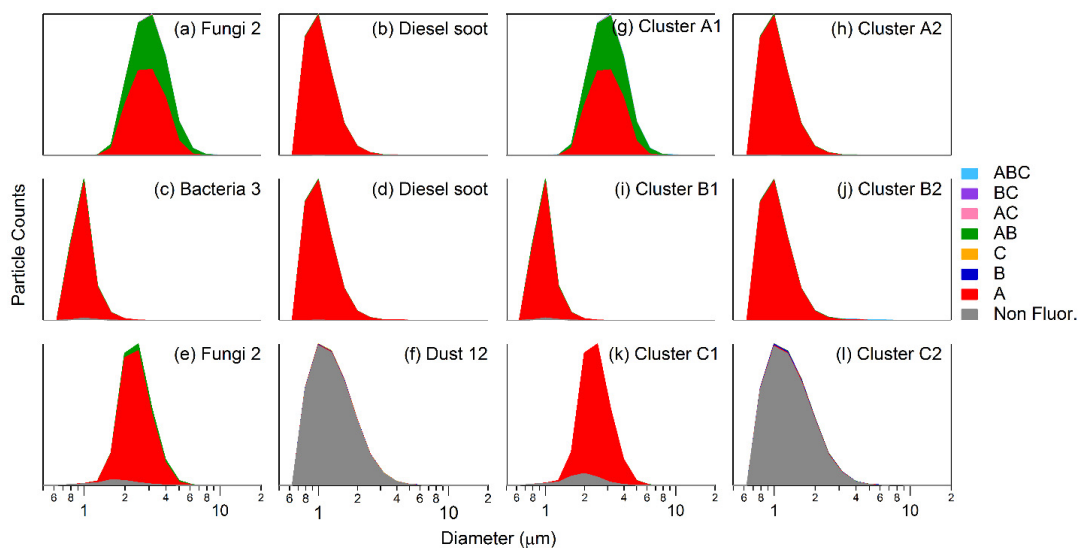|                   | A    | B   | C    | D    | E    | F    |
|-------------------|------|-----|------|------|------|------|
| **Fungi : Diesel** |      |     |      |      |      |      |
| 50:50 Ratio       | 1.1  | 0.9 | 7.2  | 4.5  | 3.6  | 0.8  |
| 80:20 Ratio       | 64.8 | 4.1 | 4.5  | 2.9  | 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 |

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

728

Figure 4. Particle type stacked category size distributions for input and output clustering results, using FT + 3σ threshold definition. Each experiment (row) shows match-ups of two particle types 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.

Atmospheric Measurement Techniques
Open Access
Discussions
EGU

**Part A: Individual Clusters**
(Particle Number)

**Part B: Grouped Clusters**
(Particle Number)

**Part C: Summary**
(Cluster Quality)

**Mixture #1: 4 Comp. - A**

| Cluster | F2 | S4 | D12 | H2 |
|---|---|---|---|---|
| 1 | 163 | 2 | 22 | 23 |
| 2 | 7 | 1 | 123 | 67 |
| 3 | 0 | 0 | 21 | 80 |
| 4 | 0 | 167 | 4 | 0 |

| Cluster | Fungi | | | Non-bio |
|---|---|---|---|---|
| 1 | 163 | | | 47 |
| 2-4 | 7 | | | 463 |

**Mixture #1**

| Total P. | Miscl. | Cat. |
|---|---|---|
| 210 | 22.4% | Fungi |
| 470 | 1.5% | Non-bio |

**Mixture #2: 4 Comp. - B**

| Cluster | F2 | S4 | D12 | WS |
|---|---|---|---|---|
| 1 | 167 | 2 | 23 | 4 |
| 2 | 2 | 3 | 88 | 10 |
| 3 | 1 | 0 | 55 | 156 |
| 4 | 0 | 165 | 4 | 0 |

| Cluster | Fungi | | | Non-bio |
|---|---|---|---|---|
| 1 | 167 | | | 29 |
| 2-4 | 3 | | | 481 |

**Mixture #2**

| Total P. | Miscl. | Cat. |
|---|---|---|
| 196 | 14.8% | Fungi |
| 484 | 0.6% | Non-bio |

**Mixture #3: High PBAP**

| Cluster | F2 | B3 | D12 | H2 | BC1 |
|---|---|---|---|---|---|
| 1 | 0 | 197 | 0 | 0 | 0 |
| 3 | 200 | 6 | 13 | 2 | 6 |
| 2 | 9 | 10 | 133 | 79 | 6 |
| 4 | 4 | 0 | 21 | 88 | 25 |
| 5 | 0 | 0 | 3 | 1 | 47 |

| Cluster | Fungi | Bacteria | Bio | Non-bio |
|---|---|---|---|---|
| 1 | 0 | 197 | | 0 |
| 3 | 200 | 6 | | 21 |
| 2,4,5 | 13 | 10 | | 403 |
| 1,3 | | | 403 | 21 |

**Mixture #3**

| Total P. | Miscl. | Cat. |
|---|---|---|
| 227 | 11.9% | Fungi |
| 197 | 0.0% | Bacteria |
| 424 | 5.0% | Bio |
| 426 | 5.4% | Non-bio |

**Mixture #4: Low PBAP**

| Cluster | F2 | B3 | S4 | D12 | H2 | BC1 | WS |
|---|---|---|---|---|---|---|---|
| 1 | 0 | 0 | 0 | 10 | 15 | 20 | 0 |
| 2 | 23 | 2 | 0 | 125 | 77 | 6 | 165 |
| 3 | 0 | 0 | 0 | 3 | 1 | 128 | 1 |
| 4 | 4 | 0 | 0 | 18 | 68 | 11 | 2 |
| 5 | 3 | 0 | 169 | 8 | 9 | 0 | 0 |
| 6 | 0 | 135 | 1 | 0 | 0 | 0 | 1 |
| 7 | 112 | 5 | 0 | 6 | 0 | 6 | 1 |

| Cluster | Fungi | Bacteria | Bio | Non-bio |
|---|---|---|---|---|
| 7 | 112 | 5 | | 13 |
| 6 | 0 | 135 | | 1 |
| 1-5 | 30 | 2 | | 836 |
| 6,7 | | | 252 | 14 |

**Mixture #4**

| Total P. | Miscl. | Cat. |
|---|---|---|
| 130 | 13.8% | Fungi |
| 136 | 0.7% | Bacteria |
| 266 | 5.3% | Bio |
| 868 | 3.7% | Non-bio |

**Mixture #5: Pollen**

| Cluster | P9 | S4 | D12 | H2 | BC1 | WS |
|---|---|---|---|---|---|---|
| 1 | 0 | 0 | 13 | 16 | 13 | 0 |
| 2 | 2 | 0 | 28 | 83 | 15 | 1 |
| 3 | 0 | 0 | 4 | 1 | 51 | 1 |
| 4 | 6 | 2 | 113 | 70 | 0 | 79 |
| 6 | 5 | 77 | 3 | 0 | 0 | 0 |
| 5 | 242 | 6 | 9 | 0 | 6 | 4 |

| Cluster | | Pollen | Non-bio |
|---|---|---|---|
| 5 | | 242 | 25 |
| 1-4,6 | | 13 | 570 |

**Mixture #5**

| Total P. | Miscl. | Cat. |
|---|---|---|
| 267 | 9.4% | Pollen |
| 583 | 2.2% | Non-bio |

**Mixture #6: Indoor Air**

| Cluster | F2 | B3 | D12 | H2 | WT |
|---|---|---|---|---|---|
| 1 | 160 | 7 | 13 | 0 | 31 |
| 4 | 0 | 154 | 0 | 0 | 0 |
| 2 | 4 | 0 | 32 | 95 | 35 |
| 3 | 6 | 9 | 125 | 75 | 62 |
| 5 | 0 | 0 | 0 | 0 | 42 |

| Cluster | Fungi | Bacteria | Bio | Non-bio |
|---|---|---|---|---|
| 1 | 160 | 7 | | 44 |
| 4 | 0 | 154 | | 0 |
| 2,3,5 | 10 | 9 | | 466 |
| 1,4 | | | 321 | 44 |

**Mixture #6**

| Total P. | Miscl. | Cat. |
|---|---|---|
| 211 | 24.2% | Fungi |
| 154 | 0.0% | Bacteria |
| 365 | 12.1% | Bio |
| 485 | 3.9% | Non-bio |

733
734
735 Figure 5. Overview of synthetic mixtures. Six mixtures shown as groups of rows, with input
736 particle fractions defined in Table 4. Part A (left columns) show particle number retrieved by
737 each individual cluster and categorized by each input particle type. Part B (middle columns)
738 show particle number categorized and grouped by particle classes (i.e. non-biological and
739 biological). Part C (right columns) show misclassification of groups of particles. Colors: light
740 green (fungal spores), blue (bacteria), pink (pollen), dark green (grouped biological), brown (all
741 non-biological).