

1 **A robust clustering algorithm for analysis of composition-dependent** 2 **organic aerosol thermal desorption measurements**

3 Ziyue Li¹, Emma L. D'Ambro^{2,3,a}, Siegfried Schobesberger^{2,4}, Cassandra J. Gaston^{2,b}, Felipe D.
4 Lopez-Hilfiker^{2,c}, Jiumeng Liu^{5,d}, John E. Shilling⁵, Joel A. Thornton^{2,3}, Christopher D. Cappa^{1,6}

5 ¹ Atmospheric Science Graduate Group, University of California, Davis, CA, USA

6 ² Department of Atmospheric Sciences, University of Washington, Seattle WA, USA

7 ³ Department of Chemistry, University of Washington, Seattle WA, USA

8 ⁴ Department of Applied Physics, University of Eastern Finland, Kuopio, Finland

9 ⁵ Atmospheric Sciences and Global Change Division, Pacific Northwest National Laboratory,
10 Richland WA, USA

11 ⁶ Department of Civil and Environmental Engineering, University of California, Davis, CA, USA

12 ^a Oak Ridge Institute for Science and Education, US Environmental Protection Agency, Research
13 Triangle Park, NC, USA

14 ^b Rosenstiel School of Marine & Atmospheric Science, University of Miami FL, USA

15 ^c TofWerk AG, Thun, Switzerland

16 ^d Now at: School of Environment, Harbin Institute of Technology, Harbin, Heilongjiang, China

17 **Abstract**

18 One of the challenges of understanding atmospheric organic aerosol (OA) stems from its complex
19 composition. Mass spectrometry is commonly used to characterize the compositional variability
20 of OA. Clustering of a mass spectral data set helps identify components that exhibit similar
21 behavior or have similar properties, facilitating understanding of sources and processes that
22 govern compositional variability. Here, we developed a clustering algorithm, Noise-Sorted
23 Scanning Clustering (NSSC), appropriate for application to thermal desorption measurements
24 from the Filter Inlet for Gases and AEROSols coupled to a chemical ionization mass spectrometer
25 (FIGAERO-CIMS). NSSC, which extends the common DBSCAN algorithm, provides a robust,
26 reproducible analysis of the FIGAERO temperature-dependent mass spectral data. The NSSC
27 allows for determination of thermal profiles for compositionally distinct clusters, increasing the
28 accessibility and enhancing the interpretation of FIGAERO data. Applications of NSSC to several
29 laboratory biogenic secondary organic aerosol (BSOA) systems demonstrate the ability of NSSC
30 to distinguish different types of thermal behaviors for the components comprising the particles
31 along with the relative mass contributions and chemical properties (e.g. average molecular
32 formula) of each cluster. For each of the systems examined, more than 80% of the total mass is
33 clustered into 9-13 clusters. Comparison of the average thermograms of the clusters between
34 systems indicate some commonality in terms of the thermal properties of different BSOA,
35 although with some system-specific behavior. Application of NSSC to sets of experiments in which
36 one experimental parameter, such as the concentration of NO, is varied demonstrates the
37 potential for clustering to elucidate the chemical factors that drive changes in the thermal
38 properties of OA. Further quantitative interpretation of the clustered thermograms followed by

39 clustering will allow for more comprehensive understanding of the thermochemical properties
40 of OA.

41 **1. Introduction**

42 Atmospheric particles are composed of hundreds to thousands of individual compounds
43 (e.g., Hamilton et al., 2004; Goldstein and Galbally, 2007), reflecting the many different sources
44 and the variety of chemical pathways that lead to their formation and growth. Various mass
45 spectrometry (MS) methods provide for characterization of this compositional variability, among
46 other techniques. Individual MS methods yield different insights into particle composition,
47 dependent upon the chemical selectivity of the method. Application of various data reduction
48 methods, such as clustering or matrix factorization, helps to reduce the inherent compositional
49 complexity and develop understanding of the sources and chemical transformations that
50 determine particle composition. Clustering and matrix factorization are complementary methods.
51 In this work, we develop and apply a new clustering method to measurements of the evolved gas
52 composition derived from thermal desorption of organic aerosol, specifically to measurements
53 from the Filter Inlet for Gases and AEROsols (Lopez-Hilfiker et al., 2014) coupled with chemical
54 ionization mass spectrometry (Lee et al., 2014) (FIGAERO-CIMS). The clustering method
55 developed here facilitates interpretation of variability in organic aerosol composition and
56 volatility, and how these depend on formation conditions.

57 Clustering methods applied across many research fields have aided in the interpretation
58 and understanding of large data sets. Clustering methods work by classifying data into several
59 groups according to the similarity between one or more properties. In the field of atmospheric
60 chemistry, clustering methods have been applied to a variety of data types. Examples include:
61 back trajectories of trace gases (Cape et al., 2000) or particles (Abdalmogith and Harrison, 2005;
62 Pinero-Garcia et al., 2015), helping to elucidate the origin and transport of pollutants; particle
63 size distributions, providing information on aerosol emission and formation (Beddows et al., 2009;
64 Wegner et al., 2012); and, the morphology of and organic functional groups comprising individual
65 particles, allowing for classification of the types of organic carbon (Takahama et al., 2007).

66 Beyond the above examples, clustering methods have been extensively applied to the
67 interpretation of single particle mass spectra, serving to characterize variability in their chemical

68 composition and identify the sources and extent of chemical processing (e.g., Gaston et al., 2013;
69 Lee et al., 2015). While clustering is a general method, a variety of specific algorithms have been
70 developed for application to a given particle mass spectral dataset. The algorithms applied to
71 analysis of single particle mass spectra include: *K*-means (Giorio et al., 2012; Liu et al., 2013; Lee
72 et al., 2015); fuzzy *c*-means (Kirchner et al., 2003; Roth et al., 2016); density-based special
73 clustering of applications with noise (DBSCAN) (Zhou et al., 2006); neural network-based
74 methods, such as an algorithm derived from Adaptive Resonance Theory (ART-2a) (Song et al.,
75 1999; Zhao et al., 2008; Giorio et al., 2012); hierarchical clustering (Murphy et al., 2003; Rebotier
76 and Prather, 2007); and, some combined algorithms (Zhao et al., 2008; Reitz et al., 2016). Each
77 clustering algorithm has strengths and weaknesses. In some cases, different algorithms are
78 equally effective and lead to similar categorization of the same data set, while in other cases
79 quite different results are obtained (Zhao et al., 2008). For example, *K*-means and ART-2a gave
80 broadly similar results on a regional particle data set (Giorio et al., 2012), and *K*-means performed
81 as well as a variant of hierarchical clustering method on four particle data sets (Rebotier and
82 Prather, 2007).

83 Here, we describe and apply a new clustering method, a novel extension of DBSCAN
84 appropriate for analysis of combined thermal desorption-mass spectral measurements of organic
85 particle composition, specifically applied to data from the FIGAERO-CIMS. FIGAERO-CIMS has
86 been increasingly used in field (e.g. Gaston et al., 2016; Lee et al., 2016; Lopez-Hilfiker et al., 2016;
87 Mohr et al., 2017; Huang et al., 2018; Le Breton et al., 2019) and laboratory studies (e.g. Lopez-
88 Hilfiker et al., 2015; D'Ambro et al., 2017; Wang and Ruiz, 2018) to develop understanding of the
89 molecular composition of organic aerosols. A key feature of FIGAERO-CIMS is the ability to
90 characterize the thermal behavior of organic compounds in particles on a near molecular level
91 (Lopez-Hilfiker et al., 2014). The use of chemical ionization, a relatively soft ionization method,
92 facilitates detection and characterization of both monomeric and oligomeric parent compounds
93 in organic aerosols. In FIGAERO-CIMS, particles are collected and then thermally desorbed, with
94 mass spectra of the evolved gases measured as a function of temperature. This can also be
95 displayed as a thermogram: the concentration of an ion or sum of ions as a function of desorption
96 temperature. The temperature at which a thermogram reaches maximum signal, or T_{max} , provide

97 information on the volatility, while particularly broad desorption shapes can indicate thermal
98 decomposition, suggesting the presence of lower volatility, possibly oligomeric, material (Lopez-
99 Hilfiker et al., 2014). A typical FIGAERO-CIMS mass spectrum of either ambient or
100 laboratory-generated organic aerosol consists of hundreds of individual ions and thermograms,
101 (D'Ambro et al., 2018; Lee et al., 2018).

102 Previous studies using FIGAERO-CIMS provided insights into particle composition, including
103 the presence of lower volatility material, based on analysis of the thermograms of several major
104 ions (Lopez-Hilfiker et al., 2014; D'Ambro et al., 2017; D'Ambro et al., 2018; Lee et al., 2018). We
105 expand on this previous work through the application of cluster analysis to FIGAERO-CIMS
106 thermograms. Clustering of FIGAERO-CIMS data provides a means to expand the understanding
107 developed from single-ion thermograms and establish the contributions of different types of
108 thermograms to the bulk particles. One previous study clustered FIGAERO-CIMS data using the
109 K-means algorithm using two parameters: the ion molecular weight and the maximum
110 desorption temperature (Faxon et al., 2018). What distinguishes our work is that we cluster the
111 thermogram across the entire desorption period for each ion, with ions grouped according to the
112 similarity of their overall volatility distribution. We have considered the performance of various
113 clustering algorithms (including K-means), ultimately concluding that a novel variant of the
114 DBSCAN algorithm, which we develop here and name noise-sorted scanning clustering (NSSC),
115 provides robust performance and has several advantages over other existing algorithms for
116 FIGAERO-CIMS data. The NSSC algorithm is applied to several laboratory data sets of secondary
117 organic aerosol (SOA) formed from various precursors and under various conditions, some are
118 previously described (D'Ambro et al., 2018). In this work we do not aim to provide comprehensive
119 interpretation of the resulting clustered thermograms in terms of their thermo-chemical
120 properties (Schobesberger et al., 2018), only to illustrate the potential of clustering to enhance
121 interpretation of FIGAERO-CIMS and other similar data.

122 **2. Clustering Method Description**

123 Application of a given clustering algorithm to a particular data type involves a number of
124 steps. Below, we discuss the specific steps for clustering of FIGAERO-CIMS data, including a

125 description of our noise-sorted scanning clustering algorithm. A brief discussion of other
126 algorithms is also provided.

127 **2.1. Data Preprocessing**

128 **2.1.1. Exclusion of anomalous thermograms**

129 The quality of the data set should be examined prior to clustering. A typical thermogram
130 exhibits a continuous evolution to a peak, peaking during a temperature ramping period, after
131 which there is a steady decrease in signal-to-background over time during a constant-
132 temperature soaking period; the background-corrected signal at all temperatures remains above
133 zero or around zero within the uncertainties. See section 3.1 for further details of the FIGAERO-
134 CIMS. An anomalous thermogram, however, contains negative signal with large magnitude.

135 Anomalous thermograms should be excluded from the clustering to assure the quality of
136 the results, although most such thermograms do not end up clustered with other ions.
137 Anomalous thermograms are identified as follows. (i) Estimate a reference noise level (σ_{ref}) for
138 each thermogram as the standard deviation of the last 100 points (corresponding to 500 seconds)
139 of the thermogram at the end of the constant-temperature soaking period, during which the
140 signals are usually relatively constant. Use of more points incorporates times when the signals
141 were still decreasing, while use of fewer points provides a less robust estimate of the noise level.
142 (ii) Find the minimum in the thermogram and calculate the average of this and the 50 points
143 (corresponding to 250 seconds, or 100 points) before and after the minimum, A_{min} . This provides
144 for consistency with the determination of σ_{ref} (iii) Identify thermograms for which $A_{\text{min}} < -3 * |\sigma_{\text{ref}}|$
145 as anomalous and exclude these associated ions from further analysis. In other words, when a
146 thermogram has a valley with averaged negative values exceeding the magnitude of three times
147 of the reference noise level, then it is considered anomalous. The specific criteria specified above
148 were determined based on consideration of thermograms from 10 distinct SOA experiments.
149 While these criteria should be robustly applicable to other FIGAERO-CIMS datasets, they can be
150 adjusted depending on the specific application, data quality, and needs.

151 Ideally, when anomalous ions are identified the original data would be inspected to identify
152 the likely origin of the anomalous behavior. Possible origins include problems with background

153 subtraction when the blank has substantially higher signal levels than the particle samples, which
154 can happen when there is residual contamination or incomplete separation of ions having the
155 same nominal mass. It is also possible that the components detected for the same ion are
156 different for the particle and blank measurements. In the example systems considered here, we
157 identified up to five anomalous ions out of what is typically a few hundred total ions.

158 In some cases, it is desirable to compare thermograms between related experiments, for
159 example the experiments discussed here that investigated the influence of NO concentration on
160 SOA formation (Section 4.3) and the impact of isothermal dilution on SOA composition and
161 volatility (Section 4.4). In such cases, ions identified as anomalous for one experiment are
162 excluded from analysis for all related experiments to ensure consistency.

163 **2.1.2. Euclidean Distance**

164 Any clustering algorithm requires a metric to determine the similarity between two
165 members in the data set. Here, we use the commonly used Euclidean Distance (ED) as the metric.
166 A smaller *ED* indicates greater similarity. A FIGAERO thermogram has *n* points, with all
167 thermograms having an equal number of points in a data set. A data set here is defined as the
168 collection of thermograms for all individual ions measured for a single desorption event. The *ED*
169 between two thermograms *a* and *b* is calculated as:

170

$$171 \quad ED_{a,b} = \sum_n \sqrt{(a_n - b_n)^2} \quad (1)$$

172

173 An individual *ED* value is obtained for every pair of ions in the mass spectrum, resulting in an *n* x
174 *n* matrix of *ED* values with the diagonal elements all zero. The signal levels between individual
175 ions differ substantially, reflecting their relative abundances. Therefore, the *ED* calculation uses
176 normalized thermograms, allowing for comparison between thermogram profiles irrespective of
177 signal magnitude. Normalization is achieved by dividing each point of the original thermogram
178 by the thermogram maximum, where the maximum is determined after smoothing using a
179 35-point boxcar moving average with the end points excluded from the smoothed thermogram.
180 Use of the smoothed maximum instead of the unsmoothed maximum reduces the influence of

181 noise on normalization. In the FIGAERO datasets used in this study, a typical thermogram has a
182 temperature resolution of $\Delta T \sim 0.7$ °C during the ramping period, and a 35-point smooth
183 corresponds to smoothing over ~ 24.5 °C. Typical FIGAERO thermograms exhibit peaks ca. 40 °C
184 wide, and thus a 35-point smoothing retains the main peak shape while reducing the influence
185 of noise. In the constant temperature part of the thermogram (soaking period), signal levels
186 change slowly with time, on average less than 5 % for a 35 points (~ 3 minutes) period, so a
187 35-point smoothing is also appropriate. We note that the unsmoothed profiles are those that are
188 normalized; smoothing relates only to determining the maximum signal values used for
189 normalization.

190 The *ED* calculation from Eqn. 1 gives equal weight to all points in the thermogram. However,
191 in a FIGAERO thermogram, equal weighting may not be appropriate. The desorption process has
192 two stages, ramping and soaking, with the soaking period comprising approximately 70% of the
193 time points in thermograms. However, most thermograms are featureless in the soaking period.
194 In contrast, many thermograms exhibit a peak, or some otherwise characteristic behavior, in the
195 ramping period. Since the behavior in the ramping period provides greater information as to the
196 overall similarity between individual thermograms, we recommend down-weighting the soaking
197 period such that the ramping and soaking periods ultimately carry approximately 4:1 weight in
198 the calculation of the *ED*. We have tested weighting of 1:1, 2:1 and 10:1. Weighting of 4:1
199 provides for the most robust clustering results for the example datasets. We do not recommend
200 completely excluding the soaking period as this period still carries informational content
201 (Schobesberger et al., 2018). Specifically, in calculating *ED* we use all data from the ramping
202 period while down-weighting the data in the soaking period by calculating and using ten-point
203 averages.

204 In summary, we calculate the *ED* based on the following steps: (i) smooth the original
205 thermogram (with absolute signal) to find the maximum value; (ii) normalize the original
206 thermogram to the smoothed maximum; (iii) average every 10 points in the soaking period; and
207 (iv) calculate the *ED* between every two normalized, down-weighted thermograms.

208 **2.1.3. Dealing with noise**

209 Noise is an inherent property of any measurement. Noise in the FIGAERO thermograms
210 results from various sources, including detector noise, background subtraction, and imperfect
211 fitting of mass spectra. Noise influences the ED calculated between two thermograms, typically
212 increasing the ED. Here, the level of noise, ξ , is characterized for each thermogram by calculating
213 the average difference between the smoothed and unsmoothed normalized thermograms for
214 the ramping period. The use of only the ramping period in assessing the noise level is consistent
215 with the generally more characteristic behavior compared to the soaking period. The use of the
216 normalized thermograms, rather than absolute, allows for comparison of noise between
217 thermograms.

218 The noise level generally varies inversely with the fractional mass contribution of the ions,
219 illustrated for a case study of the α -pinene + OH SOA (Experiment 1 in **Table 1** and **Figure 1**). This
220 indicates that ions contributing more to the total signal generally have a lower noise level.
221 Detector noise is nominally independent of ion identity, and thus the low-signal ions have
222 enhanced ξ after normalization.

223 Discussed further in section 2.3, clustering algorithms often perform poorly when overly
224 noisy data are included in the clustering. This is especially the case for algorithms such as k-means
225 and partitioning around medoids, which assign all the members to a cluster. Clustering methods
226 that do not require assignment of all members, such as DBSCAN or our NSSC, are generally less
227 sensitive to the influence of overly noisy members. However, we have found that the explicit
228 exclusion of noisy thermograms up front serves to provide for more robust behavior and also
229 removes the need to consider each noisy thermogram as a possible single-member cluster. The
230 inclusion of overly noisy peaks might obscure the underlying structure of clustered thermograms.
231 Noisy thermograms are identified as follows. First, the 5% of ions having the lowest noise are
232 identified. The ξ value of the noisiest ion from this subset of low-noise ions is defined as the
233 reference noise level, ξ_{ref} . Small differences in the choice of this threshold (e.g. using the lowest
234 7% of ions) do not materially influence the results. Ions for which $\xi_n > 3 \cdot \xi_{\text{ref}}$ are considered noisy
235 and excluded from the initial clustering. For the experiments we examined, there are 88-120 out
236 of ~300 ions left after noise screening, contributing 83.5% - 92.5% to the total particle mass.

237 2.2. Noise-sorted Scanning Clustering (NSSC)

238 2.2.1. Algorithm description

239 The noise-sorted scanning clustering (NSSC) algorithm developed here is a variant of the
240 commonly used DBSCAN. In NSSC, identification and clustering of thermograms occurs based on
241 their similarity to seed thermograms. When the ED between a given thermogram and the seed is
242 less than a specified ED criterion (ε) the two members belong to the same cluster. Importantly,
243 in NSSC the selection of the seed thermograms occurs based on their respective noise levels. The
244 least noisy thermogram is selected as the initial seed, the next noisiest is selected as the second
245 seed (assuming it is not already clustered), and so on. We have found that low-noise
246 thermograms typically have more well-defined and characteristic shapes and comprise a
247 substantial fraction of the total mass. The choice to select seeds based on the noise level leads
248 to overall more robust and reproducible clustering compared to random selection of seeds.

249 The optimal value of the distance criterion, ε , is not known *a priori*, but must be determined
250 by the user, discussed in Section 2.2.3. A valid cluster must contain at least N_{min} members,
251 inclusive of the seed. We use $N_{min} = 2$. Consideration and inspection of individual unclustered
252 thermograms exhibiting unique behavior occurs as a post-clustering process (Section 2.2.2).

253 The flow of the noise-sorted scanning clustering algorithm is shown in **Figure 2** and
254 summarized here. Clustering proceeds in two rounds. For the initial round, the thermograms are
255 sorted by the noise (ξ), and the ED values between all pairs of thermograms are calculated
256 accordingly. All of the thermograms are identified according to whether they have been already
257 used as seeds ($SEED = 0$ or 1 , with 1 for thermograms used as seeds) and whether they have been
258 already included in a cluster ($CLUSTER = 0$ or 1 , with 1 for already clustered thermograms). At the
259 start, $SEED = 0$ and $CLUSTER = 0$ for all thermograms. Clustering begins using the least noisy
260 thermogram having $SEED = 0$ and $CLUSTER = 0$ as the initial seed. The state of that seed is then
261 changed to $SEED = 1$. All thermograms having $ED < \varepsilon$ for that seed and with $CLUSTER = 0$ are
262 identified from the ED matrix; these thermograms are considered neighbors of the seed
263 thermogram. The seed does not evolve as neighbors are added to the cluster during this step. If
264 the number of neighbors plus the seed is greater than or equals N_{min} , the cluster is valid and

265 stored, with the states of all the thermograms in the cluster changed to CLUSTER = 1. Otherwise,
266 the cluster is dismissed, and CLUSTER = 0 for all the members. In this case, the current seed (with
267 SEED = 1 and CLUSTER = 0) will no longer be used as a seed in the future steps but can still end
268 up clustered as a neighbor in the other clusters. The above steps are repeated until all the
269 thermograms have either SEED = 1 or CLUSTER = 1.

270 Because a cluster must have at least N_{\min} elements, not all the thermograms may end up
271 clustered. Some of these unclustered thermograms may nonetheless have very similar shapes to
272 the clustered thermograms. Here, an iterative, second round of clustering potentially adds these
273 initially unclustered thermograms to the initial clusters, using the signal-weighted average
274 thermograms for the clusters from the first round as the initial seeds. A matrix of ED values is
275 calculated between the individual unclustered thermograms and the new seeds. For each
276 unclustered thermogram, the minimum ED , corresponding to only one of the seeds, is identified.
277 When this minimum ED is less than ε , the unclustered thermogram is added into that cluster. A
278 new signal-weighted average thermogram for the cluster is calculated and this process repeats
279 until no additional unclustered thermograms can be added to existing clusters. The mass
280 contribution of the remaining unique unclustered thermograms after this second round can be
281 substantial or negligible, ranging from <0.05% to 2.6% in the experiments presented here, and
282 depends largely on the choice of ε . Some of these unclustered thermograms are defined as
283 additional one-member clusters, discussed in the following section.

284 **2.2.2. Post-clustering Processes**

285 After thermograms are clustered, we perform two post-clustering analyses to better
286 understand the whole data set: 1) identifying additional one-member clusters and 2) sorting of
287 the clusters.

288 Some of the remaining unclustered thermograms have significant individual mass
289 contributions and should be considered as one-member clusters. The criterion of “significant”
290 mass contribution is user-defined. We recommend determining the significance criterion as
291 follows: (i) sorting all the ions (before the noise-filtering process) from largest to smallest
292 individual mass concentration; (ii) calculating the cumulative mass fraction for this sorted list;

293 and (iii) defining as “significant” all those ions contributing to a cumulative mass contribution up
294 to 80%.

295 The number of significant ions in a data set depends on the specific chemical system,
296 varying from only a few to tens of ions. Significant unclustered ions are identified as additional
297 one-member clusters. In some cases, the thermograms for these one-member clusters are
298 unique compared to the previously identified clusters. In others, their shapes are visually similar
299 to the previously identified clusters but where the one-member clusters are sufficiently distinct
300 that they were not clustered. For the purpose of automation, these one-member clusters are all
301 included in the final clustering results and the number of one-member clusters serves as one of
302 the parameters to determine the optimal ε . User can also choose to exclude them or some of
303 them manually from the final clustering results based on their judgement. For the example
304 systems considered in Section 4, there are only a few one-member clusters (ranging from 0 to 4),
305 if any, for the optimal ε used.

306 Sorting of clustered thermograms facilitates visual presentation and identification of the
307 similarities and dissimilarities among the clusters. The specific method of sorting can be varied
308 depending on the application and system under consideration. Here, we use the temperature
309 where 50% of the mass is desorbed (T_{m50}) for the weighted-average cluster thermogram as a first
310 criterion. The T_{m50} is typically similar to, but slightly larger than the temperature at which the
311 signal reaches a maximum. As such, the T_{m50} is approximately related to the saturation vapor
312 pressure of the desorbing compound, at least for compounds that desorb directly (e.g., Lopez-
313 Hilfiker et al., 2014). When two or more clustered average thermograms have identical T_{m50} , a
314 rare but occasional occurrence, they are further sorted by T_{m75} , the temperature where 75% of
315 the mass is desorbed. The temperature difference between T_{m50} and T_{m75} indicates the slope of
316 the thermogram between these two temperatures, with larger values indicating slower decay.
317 Therefore, these two parameters generally illustrate the shape of a thermogram. The T_{m50} and
318 T_{m75} are determined by calculating the cumulative desorbed mass and finding the temperatures
319 where 50% and 75% are reached.

320 The sorting process tends to organize the cluster-specific thermograms such that clusters
321 having lower peak temperatures (lower T_{m50}) and steeper downslopes after the peak (lower T_{m75})

322 come first. Thermograms of this type are indicative of major contributions from higher-volatility
323 monomers (Schobesberger et al., 2018). Thermograms having higher T_{m50} generally have broader
324 peaks, and shallower downslopes, indicative of substantial contributions from low-volatility
325 compounds or decomposition of oligomers. Further discussion of the interpretation of
326 thermogram shapes is provided in Section 3.2.

327 2.2.3. Choosing the optimal ε

328 NSSC is a distance-based clustering method, so the choice of the distance criterion, ε , is a
329 crucial step. For small ε , members within a cluster have high similarity, but few thermograms end
330 up clustered. In contrast, for large ε the majority of the thermograms are clustered into only a
331 few clusters having comparably low intra-cluster similarity. The choice of the optimal ε value is
332 guided here by consideration of several parameters that vary with ε . The overall aim is to
333 simultaneously (i) minimize the unclustered mass fraction ($f_{m,unclustered}$) while (ii) maximizing the
334 number of clusters (N_c) having two or more members and (iii) minimizing the number of one-
335 member clusters ($N_{c,one}$) yet (iv) maintain inter-cluster separation ($R_{interClst}$).

336 In general, N_c increases with ε for small ε because more thermograms of different shapes
337 get clustered and fewer thermograms remain unclustered. As ε further increases, some clusters
338 are combined and a greater number of thermograms are assigned to a single cluster.
339 Consequently, as ε increases the N_c generally increases, reaches a maximum level, and then
340 decreases. The maximum N_c and the ε at which the maximum occurs depends on the exact size
341 and the properties of dataset being examined. We have found that a typical SOA system usually
342 has 9-13 distinct thermogram clusters. We recommend selecting an ε that provides for N_c at or
343 near the maximum as this captures the greatest number of thermogram types.

344 The mass fraction of unclustered thermograms, $f_{m,unclustered}$, includes only the unclustered
345 thermograms that were not excluded based on the noise filtering. In general, a smaller $f_{m,unclustered}$
346 is preferable as this indicates a greater amount of the OA mass is included in a cluster (including
347 one-member clusters). The $f_{m,unclustered}$ generally decreases with ε , then plateaus above a certain
348 value of ε ; ideally this plateau occurs at $f_{m,unclustered} = 0$. The ε where the plateau starts is indicated
349 as ε_{MF} , where MF stands for mass fraction. Given that significant one-member clusters are

350 allowed, the unclustered thermograms that remain above ε_{MF} have individually small mass
 351 contributions and are either truly unique in their shapes or have a sufficiently high noise level
 352 that they cannot be clustered, even after the noise-screening process. We generally recommend
 353 selecting $\varepsilon \geq \varepsilon_{MF}$ to minimize the unclustered mass.

354 The number of one-member clusters, $N_{c,one}$, generally decreases with ε , as these ions are
 355 incorporated into multi-member clusters. Ideally, these one-member clusters would exhibit clear,
 356 visually distinct behavior compared to other one-member clusters and to multi-member clusters.
 357 However, we find this is often not the case, especially at smaller ε . Thus, the number of one-
 358 member clusters should generally be minimized; we suggest $N_{c,one}$ be held to five or fewer in
 359 general.

360 The inter-cluster separation parameter, $R_{interClst}$, characterizes the dissimilarity between
 361 clusters, and is the ratio between the average inter-cluster distance ($ED_{seed,avg}$) and ε , where:

$$362 \quad R_{interClst} = \frac{ED_{seed,avg}}{\varepsilon} = \frac{\sum_{i=1}^{N_{c,total}} \sum_{j=1}^{N_{c,total}} ED_{seed,i,j}}{N_{c,total} \cdot (N_{c,total} - 1) \cdot \varepsilon} \quad (2)$$

364 and $ED_{seed,i,j}$ is the distance between the seeds for the different clusters i and j and $N_{c,total} = N_c +$
 365 $N_{c,one}$. For a 2D data set, the seed can be visualized as the center of a circle and ε the radius of
 366 the circle. Thus, when $ED_{seed,i,j}/\varepsilon < 2$, the two circles defining the boundaries of these two clusters
 367 have overlapping areas. Good separation (i.e. cluster dissimilarity) is indicated when $ED_{seed,i,j}/\varepsilon >$
 368 2. Although our data set is more than two dimensions, this illustrates the idea of establishing the
 369 level of similarity (or dissimilarity) between clusters, i.e., the extent to which they are unique. We
 370 recommend selecting an ε that results in $R_{interClst} \geq 2$, when possible.

371 All four parameters should be considered when determining the optimal ε . Consideration
 372 of the parameters individually may not result in the same optimal ε . Ultimately, the user must
 373 consider each parameter and aim to select an optimal ε that balances the different information
 374 provided in each parameter. This can be achieved by plotting the above parameters as a function
 375 of ε , and then selecting as the optimal value the ε that results in (i) a small $f_{m,unclustered}$ with (ii) N_c
 376 near the maximum and (iii) a small $N_{c,one}$ and (iv) $R_{interClst}$ near or above two. In addition, visual
 377

378 comparison of the clustering results, illustrated as the average thermogram of each cluster, can
379 be helpful. For the example data considered below, we find that the optimal ϵ tends to fall within
380 a relatively narrow range of values.

381 **2.3. Alternative Clustering Methods**

382 We have alternatively considered the performance of some of the most commonly used
383 clustering algorithms (k-means, k-medoids, mean-shift, DBSCAN) and a less-commonly used one
384 (FPClustering (Gonzalez, 1985)) for interpreting FIGAERO-CIMS observations. The clustering
385 methods considered are summarized in **Table 2**, with some of their pros and cons listed, and
386 described in further detail in Appendix A. We discuss them briefly here in the context of FIGAERO-
387 CIMS data. All the methods considered require input of at least one key user-specified parameter.
388 These parameters and the associated clustering algorithms can be generally classified into two
389 categories: number-based and distance-based. Number-based clustering algorithms require
390 specifying the desired number of retrieved clusters; this includes k-means and k-medoids.
391 Number-based algorithms usually assign all members to clusters. The extent of similarity among
392 members of a cluster can vary greatly since there is no strict distance criterion for each cluster.
393 When applied to FIGAERO-CIMS thermograms, we have found these number-based algorithms
394 are particularly sensitive to the presence of noisy members and the initialization method. In
395 contrast, some clustering algorithms require specification of distance (similarity) criterion. This
396 includes the mean-shift, DBSCAN, and our NSSC algorithms. These distance-based algorithms
397 need not cluster all members of the initial population and generally emphasize intra-cluster
398 similarity or the density of the points. The methods differ in terms of the method used for
399 selection of the initial seed or center and the extent to which they emphasize point density versus
400 cluster similarity. Noisy members tend to naturally be excluded from any clusters. NSSC is a
401 variant of DBSCAN. It does, however, differ from the standard DBSCAN algorithm because NSSC
402 only searches for neighbors of the seed, while DBSCAN also searches for neighbors of the
403 neighbors. As such, the sorting of seeds by noise levels is a key aspect of the NSSC algorithm
404 which we have found provides for more robust clustering results.

405 Most of these clustering algorithms, including k-means, k-medoids, and mean-shift, are
406 initialized with a random choice of the initial cluster centers (or seeds). For large data sets, this

407 randomness usually leads to different results of clustering with different runs. The extent to
408 which this impacts analysis and clustering of FIGAERO-CIMS data is considered using SOA from
409 the α -pinene + OH SOA system as the case study (Section 4.1). For the FIGAERO-CIMS data we
410 find that the various clustering results exhibit a moderate sensitivity to how the initial seeds are
411 selected for all of these algorithms, although the final clusters are generally similar between
412 different runs for the same input parameter. This may reflect either the relatively small size of
413 the data set (~300 members originally and ~100 members after noise screening) or that there are
414 generally characteristic peak shapes with overall good separation. However, some differences
415 between independent clustering runs result, which is undesirable. For FIGAERO-CIMS data we
416 know that not all thermograms are of equal quality, i.e. they have different noise levels reflecting
417 in part their different overall contributions to the total mass. The standard clustering methods
418 do not account for this information. The NSSC algorithm developed here takes into account this
419 measure of data quality and uses it to identify the seeds for clustering. This provides for an
420 entirely reproducible clustering and generally emphasizes the behavior of the ions that
421 contribute most to the FIGAERO-CIMS signal while still allowing for consideration of contributions
422 of low-signal ions.

423 We find that different clustering algorithms can result in similar numbers of clusters with
424 the cluster-averaged thermograms having visually similar shapes when each is run with
425 appropriate user-selected parameters, although the details and robustness of each cluster vary
426 method by method. The “appropriate” parameters however are different from the “optimal”
427 parameters. There is usually different guidance for different algorithms on how to find the
428 optimal parameters that result in the greatest similarity within clusters and dissimilarity among
429 clusters. In the case of k-medoids, for example, the average silhouette indicates an optimal
430 number of clusters of two for the case study system. Yet, this is certainly too few clusters based
431 on the other methods.

432 In summary, we propose NSSC as the preferred algorithm in dealing with the FIGAERO data
433 set based on: (i) the ability to generate similar results as the other commonly used clustering
434 algorithms; (ii) good reproducibility and stability of results due to accounting for the noise of
435 individual thermograms; (iii) good control over the similarity within the clusters by using a

436 user-definable distance criterion; and (iv) a capability to identify unique thermograms as
437 one-member clusters.

438 **3. FIGAERO Measurements and Experiments**

439 **3.1. Instrument and experiment description**

440 The FIGAERO-CIMS instrument has been described previously in detail (Lee et al., 2014;
441 Lopez-Hilfiker et al., 2014). A brief description is provided here, with some additional details in
442 the Supplemental Material. The FIGAERO-CIMS measures the evolved gases from filter-collected
443 particles during temperature programmed thermal desorption. Thermal desorption of particles
444 occurs in two-stages: a “ramping” and “soaking” period. During ramping, the temperature
445 increases from room temperature to 200 °C, typically at 10 °C min⁻¹. Most OA mass desorbs
446 during the ramping stage. The temperature is held at 200 °C for ca. 30–40 mins during the soaking
447 period to facilitate evaporation of the remaining, low-volatility organic mass from the filter. The
448 evolved gas-phase compounds are measured using CIMS with the iodide (I⁻) reagent ion,
449 appropriate for characterization of generally highly oxygenated components comprising most
450 secondary organic aerosol (Lopez-Hilfiker et al., 2016; Isaacman-VanWertz et al., 2017; Lee et al.,
451 2018). The resulting signal or mass concentration versus temperature (or equivalently time)
452 curves for each ion constitute a thermogram. All individual thermograms are background
453 corrected by subtracting the observed thermograms from appropriate blank experiments. The
454 overall bulk thermogram is obtained by summing together the individual thermograms.

455 Several example applications of the clustering on FIGAERO-CIMS data are discussed in
456 Section 4. These cover laboratory experiments on SOA derived from: (1) OH + α -pinene and (2)
457 OH + Δ -3-carene, both at low-NO_x conditions; (3) OH + α -pinene as a function of [NO]; and (4)
458 O₃ + α -pinene, but where the SOA is allowed to isothermally evaporate at 80% RH for varying
459 amounts of time prior to thermal desorption. These experiments are summarized in **Table 1**, with
460 further details in the Supplemental Material and associated publications (D'Ambro et al., 2018;
461 D'Ambro et al., 2019); all data are publicly available (Cappa et al., 2019). All the experiments were
462 done in a 10.6 m³ Teflon environmental chamber at Pacific Northwest National Laboratory (PNNL)
463 (Liu et al., 2012; Liu et al., 2016).

3.2. General interpretation of FIGAERO-CIMS thermograms

This work focuses on development of the clustering method, rather than on interpretation of the FIGAERO-CIMS thermograms; an illustrative thermogram is shown in **Figure 3b**. However, discussion of the clustering results is aided by a general understanding of how FIGAERO-CIMS thermograms have been previously interpreted. Ions contributed by semi- and low-volatility compounds that desorb directly tend to exhibit strongly peaked, Gaussian-like thermograms with single-mode peaks between around 50 °C to 120 °C; the lower the peak desorption temperature (T_{peak}) the higher the volatility of the desorbing compound (Lopez-Hilfiker et al., 2014; 2015). We therefore refer to thermograms, or portions of thermograms, having this general shape as the “monomeric” content of the ion hereafter; direct evaporation of thermally stable dimers or other oligomers is possible, although will typically occur at higher temperatures due to the comparably lower volatility of these compounds. When multiple monomeric compounds having different vapor pressures contribute to the same ion, the resulting thermogram exhibits a broader peak and shallower slopes or, in particular cases, multiple, distinct peaks (Lopez-Hilfiker et al., 2015). However, very broad thermograms, especially those that peak at higher temperatures (> 120 °C or so), can also indicate contributions from thermal decomposition of very low-volatility monomers, dimers, and oligomers (Lopez-Hilfiker et al., 2015; Gaston et al., 2016; Schobesberger et al., 2018). Dimers and oligomers can evaporate directly, without thermal decomposition, as observed for isoprene-derived SOA (D'Ambro et al., 2017) and ambient monoterpene oxidation products (Mohr et al., 2017). However, fragments of dimers or oligomers are generally more abundant, indicating the importance of thermal decomposition for desorption of these low-volatility compounds. Both direct evaporation of extremely low-volatility compounds and decomposition of large molecules or oligomers can lead to high signal levels above ~120 °C. We refer to both peaks and the slowly varying signal above ~120 °C as the “oligomeric” content of the ion hereafter. We use the terms monomer and oligomer in a qualitative manner. A more quantitative analysis of the thermograms can help distinguish between direct evaporation, thermal decomposition, and the contributions of monomers versus oligomers (Schobesberger et al., 2018), yet is beyond the scope of the current work.

492 4. Example Applications

493 To illustrate the broad utility of NSSC for interpretation and analysis of FIGAERO-CIMS data,
494 we apply NSSC to the laboratory-generated SOA systems described above. The systems include:
495 SOA formed from a single precursor under NO_x-free conditions; SOA formed from a single
496 precursor as a function of input [NO]; and, SOA formed from a single precursor with thermal
497 desorption following isothermal evaporation.

498 4.1. α -pinene + OH SOA

499 A total of 298 ions were characterized by FIGAERO-CIMS for SOA generated from the
500 α -pinene + OH reaction (**Table 1**). Four ions were characterized as anomalous and excluded from
501 further analysis (see Section 2.1.1). The mass concentration of each ion was calculated by
502 integrating the signal across the entire desorption period and assuming an equal sensitivity of
503 CIMS for all the compounds. The total mass concentration is the sum of all the non-anomalous
504 ions. The mass spectrum and bulk thermogram of the remaining 294 ions are shown in **Figure 3**,
505 with the bulk thermogram shown versus both temperature (**Figure 3b**) and time (**Figure 3c**) to
506 illustrate the difference between the ramping and soaking periods. The individual thermograms
507 exhibited a variety of shapes. The noise threshold for this data set was $\xi_{\text{ref}} = 0.020893$. A total of
508 188 ions were screened out via noise filtering. The remaining 106 ions contribute 92.5% to the
509 total mass detected by FIGAERO-CIMS. The optimal ε was established through consideration of
510 the co-dependencies of N_c , $N_{c,\text{total}}$, $f_{m,\text{unclustered}}$ and $R_{\text{interClst}}$ on ε (**Figure 4; Table 3**). For this data
511 set, we determine the optimal $\varepsilon = 2.6$. Choice of a much smaller ε , around 1.5, gives a maximum
512 in N_c , but leaves a large fraction of the mass unclustered. Choice of $\varepsilon = 2.1$ or 2.2 yields larger N_c
513 and $R_{\text{interClst}}$ than $\varepsilon = 2.6$, with a reasonably small $f_{m,\text{unclustered}}$. However, there is one type of
514 thermogram (Clst#11 in **Figure 5**) that is only captured with $\varepsilon \geq 2.6$ and this yields $f_{m,\text{unclustered}} = 0$.
515 Using $\varepsilon \geq 2.7$ also yields $f_{m,\text{unclustered}} = 0$ and $N_{c,\text{one}} = 0$, but N_c and $R_{\text{interClst}}$ decrease from $\varepsilon = 2.6$,
516 indicating increasing similarity between clusters with fewer types of shapes captured. The choice
517 of $\varepsilon = 2.6$ provides a compromise between maximizing N_c , minimizing $f_{m,\text{unclustered}}$, and keeping
518 $R_{\text{interClst}}$ above two. The parameters and thresholds used for this data set are summarized in **Table**
519 **3**.

520 A total of 11 clusters are identified with no one-member clusters. The unweighted and
521 mass-weighted average thermograms for each cluster are shown along with the thermograms of
522 individual members in **Figure 5a**. The differences between weighted and unweighted average
523 clusters are negligible, in general. Clusters are organized and numbered (as Clst#*N*) from low to
524 high T_{m50} , with deeper to shallower downslope. Clst#1 through Clst#6 all have a clear peak below
525 120 °C, but with different peak widths and downslopes. Clst#7 and Clst#8 are a bit noisier with
526 only a few members each, exhibiting a sharp upslope and shallow downslope. Clst#9 has a very
527 broad peak. Clst#10 peaks at around 150 °C after an initial rise and temporary plateau. Clst#11
528 exhibits behavior somewhat like Clst#10, but with a peak that occurs just into the soaking period,
529 evident if viewed in time space, at 200 °C with a rapid drop afterwards.

530 The total mass concentration of a given cluster ($M_{c,N}$) is the sum across all cluster members,
531 calculated by integrating the summed mass concentration across the entire desorption period.
532 The percentage mass contribution of each cluster, and of the unclustered and the noise-filtered
533 ions, as well as the number of members for each cluster are shown in **Figure 5b** and **Table S1**.
534 Clst#2 and Clst#3 contain the majority of the mass (20.1% and 44.3%, respectively) and consist
535 of nearly half of the clustered ions (11 and 42, respectively). Clst#4 and Clst#9 also contain a
536 notable percentage of the total mass (8.2% and 9.8%, respectively) and include a notable number
537 of ions (13 and 17, respectively). Other clusters contribute relatively little to the total mass and
538 contain a small fraction of ions.

539 The mass-weighted average molecular formulas ($C_xH_yO_zN_m$) differ between clusters, as do
540 the O:C and H:C atomic ratios (**Table S1**). There is no clear relationship between T_{m50} (or cluster
541 number) and the number of carbon atoms, MW, or O:C. There is, however, a reasonable, inverse
542 correlation between T_{m50} and H:C ($r^2 = 0.78$). The number of carbon atoms is notably larger for
543 Cluster 6 ($x = 11.1$) and Cluster 7 ($x = 15.3$); if those two clusters are excluded there is an inverse
544 relationship between T_{m50} and the number of carbon atoms ($r^2 = 0.79$) and with MW ($r^2 = 0.59$).
545 While the reason for these two clusters having comparably large numbers of carbon atoms is
546 unknown, this nonetheless suggests that the contribution of oligomer decomposition might
547 increase for clusters having higher T_{m50} values.

548 Interpretation of previous FIGAERO-CIMS studies have largely focused on the behavior of
549 the bulk thermogram or of several major ions or sums of ions based on common factors such as
550 the number of carbon atoms (Lopez-Hilfiker et al., 2016; D'Ambro et al., 2017; D'Ambro et al.,
551 2018; Stolzenburg et al., 2018; Wang and Ruiz, 2018; Joo et al., 2019). The normalized
552 thermograms of the top five ions contributing most to the total mass for the experiments here
553 are shown in **Figure 5c**, along with the bulk thermogram. Together these five ions make up nearly
554 30% of the total mass, and exhibit very similar thermogram shapes to each other and to the bulk
555 thermogram and belong solely to either Clst#2 or Clst#3. Thus, examining these ions only would
556 capture only a fraction of the overall diversity in thermal behaviors. The clustering method
557 developed here provides a means to investigate more comprehensively the variability in volatility
558 between aerosol components.

559 **4.2. Δ -3-carene + OH SOA**

560 A total of 298 ions were characterized by FIGAERO-CIMS for SOA generated from the
561 reaction of Δ -3-carene + OH (**Table 1**). Five were identified as having anomalous thermograms
562 and excluded from further analysis. The mass spectrum and bulk thermograms of Δ -3-carene +
563 OH SOA are shown in **Figure 6**. Compared to the α -pinene +OH SOA described above, the mass
564 spectrum of Δ -3-carene SOA is quite different, with one ion ($C_8H_{12}O_5$) dominant. The bulk
565 thermograms of the two SOA systems both look bell-like, but with the Δ -3-carene SOA
566 thermogram having a peak temperature ca. 9 °C higher. After noise-filtering, 110 ions remained
567 for clustering, contributing 90.7% to the total mass. The optimal $\varepsilon = 2.1$, established again by
568 considering the system-specific dependence of N_c , $N_{c,one}$, $f_{m,unclustered}$ and $R_{interClst}$ on ε (**Figure S1**),
569 with the parameters and thresholds summarized in **Table 3**.

570 Ten clusters are identified, including one one-member cluster, with thermograms shown in
571 **Figure 7a** and the mass contribution and number of ions in a cluster in **Figure 7b**. Chemical
572 properties of each cluster are summarized in **Table S2**. The general characteristics of
573 thermograms identified in the Δ -3-carene + OH SOA are similar to those of low- NO_x α -pinene +
574 OH SOA described above, but with different mass contributions. For example, Clst#4 has nearly
575 identical shape of the thermogram as Clst#3 in the α -pinene SOA but contributes less to the total

576 mass, 28.0% compared to 44.3%. Clst#6 in the Δ -3-carene SOA contributes 14.8% to the total
577 mass and resembles Clst#5 in the α -pinene SOA, which contributes only 4.0% to the total mass.

578 In general, Clst#1 – 6 in the Δ -3-carene SOA all exhibit a peak below 120 °C, with clear peaks
579 of varying width and downslopes of varying steepness, but nominally in order of narrow to wide
580 and steep to shallow, respectively. These clusters carry the majority of the desorbed mass. Clst#7
581 and Clst#8 both exhibit relatively flat thermograms in the ramping period after their initial rise,
582 and contribute 9% to the total mass. Clst#9 has a peak temperature above 150 °C and Clst#10
583 reaches a maximum during the soaking period. These last two clusters contribute little to the
584 total mass (0.6% and 0.3%, respectively).

585 The thermograms of the five largest ions are shown in **Figure 7c**. These five ions together
586 carry ~35% of the SOA mass. A wider variety of thermogram shapes are captured by the top five
587 ions compared to the α -pinene SOA system. However, thermograms characteristic of Clst#7–10
588 are not represented by these top five ions; this remains true even if the top 10 ions are
589 considered (not shown).

590 There are ultimately three major differences between the two SOA systems. For one, there
591 is a different relationship between fractional contribution and cluster number (and thus $T_{m,50}$)
592 between the two. Secondly, the α -pinene SOA contains ions with especially narrow peaks at ca.
593 100 °C (i.e., Clst#7 & 8), that are not observed with Δ -3-carene SOA (compare **Figure 5** with **Figure**
594 **7**). Lastly, the thermograms of the top five ions for Δ -3-carene SOA differ to a greater extent than
595 for α -pinene SOA. Although we are unable to determine the reasons for these differences here,
596 this illustrates the potential for clustering to help identify and understand differences between
597 different SOA systems.

598 **4.3. α -pinene + OH + NO SOA**

599 Thermograms from SOA generated from the reaction of α -pinene + OH at varying NO
600 concentrations (5 ppb, 10 ppb and 25 ppb; **Table 1**) are considered as a set of experiments.
601 Together, differences between them illustrate the impact of changes to the fate of RO₂ peroxy
602 radical intermediates on the SOA composition and thermal properties (Praske et al., 2018; Zhao
603 et al., 2018). Clustering proceeds here using two complementary approaches. In the single

604 clustering method, clustering is performed for one reference experiment (i.e., at one NO
605 concentration, 5 ppb, Expt#3a). Then, average thermograms are calculated for the other
606 experiments in the set using the same cluster members as identified in the reference experiment.
607 In the multiple clustering method, clusters are independently determined for each experiment in
608 the set, and the shapes, relative abundances, and contributing ions are compared between
609 experiments. For all three experiments, the same initial set of 298 ions were characterized by
610 FIGAERO-CIMS.

611 **4.3.1. Single Clustering**

612 The ions identified as anomalous in each experiment differed. This most likely results from
613 shifts in the background signal levels between experiments. To maintain consistency between
614 the three experiments, ions identified as anomalous in any of the experiments were excluded
615 from all the experiments, with four ions excluded in total. A total of 88 ions were kept for
616 clustering after noise-filtering using the 5 ppb NO reference experiment, contributing 84.5% to
617 the total mass. The optimal $\varepsilon = 2.2$ (**Figure S2** and **Table 3**), resulting in ten clusters with one
618 one-member cluster. The same sets of ions were then used to calculate the cluster-average
619 thermograms for the 10 ppb and 25 ppb NO experiments. Chemical characteristics of the clusters
620 are summarized in **Table S3**.

621 Mass spectra for the three experiments are compared in **Figure 8a** and the bulk
622 thermograms shown in **Figure 8b** and c. The 5 ppb NO and 10 ppb NO SOA mass spectra are
623 nearly identical. The mass spectrum for the 25 ppb NO experiment, however, exhibits a notable
624 shift of the most abundant ions towards lower m/z . The bulk thermograms for the 5 ppb and 10
625 ppb NO experiments are nearly identical, peaking near 80 °C. The 25 ppb NO bulk thermogram
626 similarly peaks near 80 °C, but exhibits a much slower decay as temperature increases further.
627 Additionally, the change in slope at the transition from the ramping to soaking period is more
628 pronounced in the 25 ppb NO experiment. Overall, a greater fraction of the mass desorbs above
629 100 °C and during the soaking period for the 25 ppb NO experiment compared to lower-NO
630 experiments.

631 Despite the differences in the bulk thermograms, the shapes of the weighted-average
632 thermograms of clusters for all the NO experiments are generally similar, with the exception of
633 Clst#6 (**Figure 9a**). In particular, the 25 ppb thermogram shape of Clst#6 differs substantially from
634 those of low-NO conditions, with a much reduced initial peak (around 80 °C) and an more
635 pronounced second peak at high temperature (around 200 °C). However, this cluster contributes
636 negligibly to the overall mass. There is some suggestion of similar behavior for Clst#10, although
637 to a lesser extent. For the three most abundant clusters, Clst#1, 2 and 4, there is a slightly
638 increased relative contribution of the 100-200 °C tail for 25 ppb NO, consistent with differences
639 in the bulk thermograms.

640 The most notable NO-dependent change is in the relative abundances of the clusters
641 between the 5 and 10 ppb NO experiments and the 25 ppb NO experiment (**Figure 9b**). The
642 cluster mass fractions are nearly identical between the 5 and 10 ppb NO experiments. The
643 relative contributions of higher-number clusters (which have been ordered according to
644 increasing $T_{m,50}$) increase for the 25 ppb NO experiment. This is consistent with the increased
645 persistence of the 25 ppb NO bulk thermogram to higher temperatures and the nearly identical
646 nature of the 5 ppb and 10 ppb NO bulk thermograms (**Figure 8b**). The clustering analysis suggests
647 that differences in the bulk thermogram arise from shifts in the relative contributions of the
648 various SOA components that result from the altered photochemical environment. These
649 observations generally suggest an increasing fraction of oligomeric content, or less-volatile
650 compounds, formed in the particle phase—or potentially the gas phase—when the SOA was
651 generated under higher chamber NO conditions (Schobesberger et al., 2018).

652 **4.3.2. Multiple Clustering**

653 With multiple clustering, each experiment was processed and clustered independently,
654 with experiment-specific ξ_{ref} , N_c , and ε , among other parameters (**Figure S4** and **Table 3**). The
655 clustered thermograms from the three experiments are compared in **Figure 10a-c**. The number
656 of clusters identified increases with NO concentration. Comparison between the shapes of the
657 clusters from the 5 ppb NO (**Figure 10a**) and 10 ppb NO (**Figure 10b**) experiments indicates
658 generally similar types of thermograms, consistent with the single clustering method. Ten of the
659 11 total 10 ppb clusters match with a 5 ppb cluster. The one additional, unique cluster at 10 ppb

660 NO (Clst#9), is a one-member cluster with a sharp, narrow peak at low temperatures and a
661 broader, shallow second peak at high temperatures. This ion was filtered out due to high noise
662 level in the 5 ppb NO experiment.

663 The 25 ppb NO experiment (**Figure 10c**) results in more clusters compared to the lower NO
664 experiments; 13 for the 25 ppb NO experiment versus 10 and 11 for the 5 and 10 ppb experiments,
665 respectively. Some of the 25 ppb NO clusters have shapes similar to the lower NO experiments,
666 but many differ substantially. For example, two of the unique 25 ppb NO clusters (Clst#12 and
667 #13) have thermograms for which the signal increases continuously through the ramping period
668 and even into the soaking period. These clusters were not found in the single clustering analysis
669 because the 5 ppb NO experiment was used as the reference.

670 The new types of thermograms observed in the 25 ppb NO experiment indicates either
671 formation of new compounds or a change in the relative contributions of different components
672 to the same ions. Either could result from a change in the fate of the peroxy radical intermediates
673 as the NO concentration increases, leading to notably different products. There were numerous
674 nitrogen-containing ions observed for the three experiments. These N-containing ions belong to
675 Clst#1 – 7 for all the three [NO] conditions (**Table S4**). The higher-number clusters did not include
676 N-containing ions, also indicating a limited influence of the N-containing products on these lower-
677 volatility thermograms, although fragmentation complicates the interpretation. Overall, the
678 formation of new N-containing compounds at the high NO condition does not seem to explain
679 the unique thermograms in the 25 ppb NO experiments.

680 The percent contribution of different clusters to total mass, along with the noise-filtered
681 and unclustered ions, differ between experiments (**Figure 10d**). Note that for the multiple
682 clustering method, clusters having the same index number are not necessarily directly
683 comparable between experiments because different sets of ions are included. For example, while
684 Clst#1 in the 5 ppb and 10 ppb NO experiments are comparable, the most similar cluster in the
685 25 ppb experiment is Clst#2. Nonetheless, there are some common features shared by the same,
686 or closely indexed, clusters. For example, Clst#1 – 4 in all three experiments exhibit a narrow,
687 single peak with the peak temperature below 120 °C. The mass contribution of Clst#1 – 4 is similar
688 between the 5 and 10 ppb NO experiment, but ~15% lower in the 25 ppb NO experiment. Clusters

689 that reach their maximum signal at or above 150 °C (Clst#9, 10 for 5 ppb, Clst#10, 11 for 10 ppb
690 and Clst#10 – 13 for 25 ppb) together contribute ~6% in the low NO experiments and ~13% in
691 the high NO experiments. Thus, there is some evidence that at higher NO there is an increased
692 contribution of oligomeric compounds, indicated by the increased contribution of clusters that
693 peak at higher temperatures and exhibit broader overall thermograms. However, overall these
694 observations suggest complex shifts in the distribution of products, both monomeric and
695 oligomeric, with sufficient increases in NO to change the fate of the peroxy radical intermediates.

696 **4.4. α -pinene + O₃ SOA**

697 SOA formed from dark ozonolysis of α -pinene was collected and then allowed to
698 isothermally evaporate for varying amounts of time (0 h, 1 h, 3 h, 6 h and 24 h) before thermal
699 desorption (**Table 1**, Expt#4). As above for the SOA formed at varying NO concentrations, these
700 experiments are considered as a set and interpreted using both the single-clustering and
701 multiple-clustering approaches. The single-clustering approach uses the 0 h (no-wait) experiment
702 as the reference for initial clustering. In this set of experiments, 312 ions were characterized by
703 FIGAERO-CIMS for each experiment.

704 **4.4.1. Single Clustering**

705 Only a few ions, if any, were identified as anomalous in each experiment; a total of ten ions
706 were removed from all the experiments to maintain consistency between experiments. The mass
707 spectra and bulk thermograms of the remaining 302 ions for the five experiments are shown in
708 **Figure 11**. As the isothermal evaporation time increases, the mass spectrum changes significantly,
709 as previously reported by D'Ambro et al. (2018). In the no-wait experiment, the mass spectrum
710 is dominated by one ion, C₁₀H₁₄O₆. Upon isothermal evaporation, the relative abundance of this
711 ion notably decreases, with the extent of decrease increasing with wait time; over time, a greater
712 number of ions contribute to the total mass, both at lower and higher m/z . With isothermal
713 evaporation, the bulk thermograms also exhibit a shift from a more peaked shape, reminiscent
714 of that from a single compound (Lopez-Hilfiker et al., 2014), to a more flattened peak with a
715 shallower rise (**Figure 11**). In other words, with increasing isothermal evaporation the majority
716 of the mass desorbed during thermal desorption shifts from a lower to higher temperature region.

717 This behavior largely reflects the loss of comparably more volatile compounds during isothermal
718 evaporation, leaving behind SOA that is overall less volatile (**Figure S6a**). It can also in part be due
719 to higher molecular weight, lower volatility compounds being produced with time via accretion
720 reactions in the condensed phase.

721 There are 12 clusters determined from the no-wait experiment, exhibiting a wide variety of
722 the shapes (**Figure 12a**), with the parameters used for data pre-processing and clustering
723 reported in **Table 3** and shown in **Figure S5**. Focusing first on the no-wait experiment, the cluster
724 thermogram shapes include those having clear peaks at relatively low temperatures (~ 60 °C) and
725 with a sharp rise and fall (e.g., Clst#1-3), those having sharp peaks at relatively low temperatures
726 but with a shallow downward slope (e.g., Clst#6), those with a broad peak at somewhat higher
727 temperatures (~ 100 °C) and long tails (e.g., Clst#7), and those having a wide peak at even higher
728 temperatures ~ 120 °C with a very broad rise and fall (e.g., Clst#10).

729 Changes to the shapes of the thermograms that occur upon isothermal evaporation differ
730 between the clusters. Some of the clusters exhibit almost step changes from the no-wait to the
731 longer time experiments (e.g., Clst#2 and 6), while others exhibit more continuous changes (e.g.,
732 Clst#3 and 5). However, in all cases the clusters shift to have peaks that occur at higher
733 temperatures with generally broader thermograms. In other words, the T_{m50} of all the clusters
734 increase as a function of evaporation time, but with larger increases observed for the clusters
735 having initially lower $T_{m,50}$ (**Figure 12b**). For some of the clusters with a clear peak below 100 °C,
736 such as Clst#1–6, the peaks broaden to become less obvious and shift to higher temperatures
737 with longer isothermal evaporation. For clusters that originally have very wide peaks, such as
738 Clst#8–10 and 12, isothermal evaporation engenders a general shift in the thermograms towards
739 higher temperatures. Different from the clusters described above, thermograms for two clusters,
740 Clst#7 and Clst#11, exhibit only minor shift of peak temperature and shapes. Thermograms of
741 these two clusters share the common features of a moderate-width peak that reaches a
742 maximum between 100 – 120 °C. The T_{m50} of these two clusters correspondingly exhibit small
743 changes compared to other clusters.

744 Isothermal evaporation generally leads to a reduction of the monomeric character of
745 clusters, leaving behind components that exhibit increased oligomeric content. Differences in

746 how the individual cluster thermograms evolve with isothermal evaporation are therefore likely
747 indicative of differing relative contributions of monomeric versus oligomeric components. For
748 example, Clst#1 and Clst#10 have distinctly different shapes in the 0-h wait experiment, but very
749 similar shapes in the 24-h wait experiment. This indicates that ions in Clst#1 are not contributed
750 from a single component, as might be inferred from the single-mode peak in the 0-h wait
751 experiment. Instead, they are contributed by multiple components, though initially dominated
752 by monomeric compounds, so the shift in peak temperature and broadness is substantial. On the
753 other hand, ions in Clst#10 must also derive from multiple components, but with only a small
754 fraction of monomeric compounds that evaporate in the 24 hours. Consequently, the loss of
755 low-temperature mass is apparent yet small. In contrast, ions in clusters such as Clst#7 and 11
756 must be composed of only low-volatility components because they exhibit minimal changes in
757 the thermograms shapes.

758 The extent of mass loss with isothermal evaporation differs between clusters. In general,
759 clusters that exhibit larger changes in shape have greater total mass loss, although with variability
760 (**Figure S6c**). Consequently, the mass contributions of the clusters evolve with isothermal
761 evaporation (**Figure 12b**). The contribution of Clst#1 decreases significantly and most notably as
762 wait time increases. The most prominent ion in the no-wait experiment, $C_{10}H_{14}O_6$, is grouped in
763 Clst#1. The continuous mass loss of Clst#1 indicates the rapid evaporation of its members. The
764 mass contributions of the other clusters that exhibited similar changes in shape as Clst#1 (Clst#3,
765 5, and 6) remain comparably constant, although with Clst#3 decreasing slightly. The relative
766 abundances of the clusters for which the thermograms shapes changed negligibly (Clst#7 and 11)
767 increase continually, implying of the slowest evaporation of the ions in these two clusters in the
768 24-hr evaporation period.

769 For comparison, D'Ambro et al. (2018) reported changes in the shapes of the thermograms
770 for the five most abundant individual ions from the no-wait to 24-hr experiment, together
771 carrying ~15% of the particle mass. They observed the individual ion thermograms generally all
772 evolved in a manner similar to our Clst#1, 3 and 5, shifting from narrower, more peaked profiles
773 towards broader profiles with a shallower rise, less evident peak, and increased evaporation at
774 higher temperatures. Here, with the clustering of data, we are able to track the change of thermal

775 behaviors of ions carrying ~87% of the initial mass. We are able to confirm that ~70 % of the mass
776 exhibit similar thermal behaviors and responses to isothermal evaporation as the top five ions.
777 However, we are also able to identify another ~17% of the mass having initial thermograms not
778 characterized by the top five ions, including 12% of the mass (Clst#7 and 11) that behaves
779 distinctly different upon evaporation at room temperature.

780 **4.4.2. Multiple Clustering**

781 The number of clusters identified with the multiple-clustering method, using experiment-
782 specific optimal ε values (**Table 3** and **Figure S7**), decreases with isothermal evaporation time,
783 from 13 (no-wait) to 12 (1 h) to 11 (3 h) and then to 9 (6 h and 24 h) (**Figure 13b-f**). The noise
784 levels of the thermograms increase with evaporation time due to decreasing absolute particle
785 mass. Nonetheless, the typical shapes of the cluster-specific thermograms clearly evolve with
786 increasing isothermal evaporation. For short isothermal evaporation times, many cluster-specific
787 thermogram profiles are relatively narrow, peaking at lower temperatures (70-120 °C) and with
788 rapid rises and evident downslopes. For longer isothermal evaporation times, the cluster-specific
789 profiles instead have broad peaks with slow rises and most of the mass desorbing at higher
790 temperatures.

791 To aid further general interpretation, the cluster-specific thermograms with $T_{m50} < 120$ °C
792 are grouped together as higher-volatility clusters. The number of higher-volatility clusters
793 decreases with isothermal evaporation, from ten for the no-wait experiment, to five in the 1-h
794 experiment, two in the 3-h and 6-h experiment, to none in the 24-h experiment (**Figure 14**). The
795 mass contributions of the higher-volatility clusters decrease from 81.9% to 60.4%, 17.2%, 9.4%
796 and to 0.0%, with increasing isothermal evaporation time. This overall behavior is consistent with
797 results from the single-clustering method and indicates the compounds with a wide range of
798 volatilities make up much of the mass in the initial particles, while the SOA after isothermal
799 evaporation is composed of compounds having lower volatilities.

800 After isothermal evaporation, some cluster-specific thermograms have signals that increase
801 continuously during the ramping period, for example Clst#11 and 12 in the 1-h experiment; such
802 clusters were not observed in the no-wait experiment. The relative abundance of these very low-

803 volatility clusters increases with isothermal evaporation, from 1.7% in the 1-h experiment
804 (Clst#11 and 12) to 13.4% in the 24-hr experiment (Clst#7 and 9). The absence of these clusters
805 for the no-wait experiment suggests that they are formed over time through condensed-phase
806 reactions. Their increasing contribution over time may reflect both evaporation of higher
807 volatility components and continued formation. Clusters having thermograms with very broad
808 peaks, such as Clst#11 and 13 in the 0-h experiment are also observed in all the other experiments,
809 with increasing contribution to the total mass.

810 The multiple-clustering method reveals the disappearance of certain types of thermograms,
811 (e.g., the no-wait Clst#3) and the emergence of other types of thermograms (e.g., the 1-h Clst#11)
812 as evaporation time increases. This complements the single-clustering method, which illustrates
813 gradual changes in the shapes of cluster-specific thermograms, by allowing for identification of
814 completely new thermogram shapes and divergent behavior between ions within initial clusters.
815 The multiple-clustering method also confirms the decrease of the diversity of the desorption
816 profiles, as suggested by the single-clustering method. The two methods complement each other
817 and together provide a detailed look into (i) how the desorption profiles of sets of ions evolve
818 with isothermal evaporation and (ii) how the fraction of different types of thermograms change
819 with evaporation time.

820 **5. Conclusions**

821 We developed a new clustering algorithm, the noise-sorted scanning clustering (NSSC)
822 algorithm, for application to FIGAERO-CIMS data sets. The NSSC algorithm provides a robust
823 method for clustering of FIGAERO-CIMS thermograms having distinct thermal desorption profiles
824 and of determining the mass contribution of each cluster. Each of the ions contributing to a
825 cluster results from one or more molecules sharing similar thermochemical properties. These
826 molecules either evaporate directly or decompose and then evaporate. Compared to other
827 existing clustering algorithms, NSSC is strictly similarity-based, reproducible, and takes into
828 consideration differences in noise levels between individual ions. The application of NSSC has the
829 potential to make FIGAERO data more accessible to the atmospheric chemistry community.

830 For the four different SOA systems we examined, more than 80% of the total mass is
831 clustered, with the number of clusters ranging from 9 to 13. The shapes of the cluster-specific
832 average thermograms exhibit substantial variation for a given system. Some have relatively sharp
833 peaks, others broad peaks with slowly decreasing signal as heating continues, and others still
834 having signals that continually increase up to very high temperatures or long desorption times.
835 The mass contribution of a cluster varies from 0.2% to 44.3%. A few (2-3) clusters usually contain
836 more than 50% of the total mass in all the chemical systems examined. Comparison of the cluster-
837 specific thermogram shapes between different SOA systems allows for qualitative assessment of
838 the similarity or uniqueness.

839 We also demonstrated the potential of the NSSC for guiding interpretation of sets of
840 experiments where one experimental condition varies (e.g., NO concentration and evaporation
841 time). For such experiments, two complementary methods are suggested: (i) the single clustering
842 method, where one experiment is used to determine the ions belonging to individual clusters
843 and then clusters comprising the same ions are calculated for the other experiments, and (ii) the
844 multiple clustering method, where each experiment is clustered independently and then
845 compared. The first approach helps establish how the properties of individual clusters evolve as
846 a set, while the second approach helps identify changes in the diversity of cluster-specific
847 thermogram shapes, properties, and mass contributions. The two approaches complement each
848 other and provide guidance for future efforts to cluster ambient observations having long time-
849 series.

850 This paper focuses only on the description of the clustering algorithm and its potential as a
851 tool to characterize the thermal properties of organic aerosol in further detail. [The application of
852 NSSC can be potentially expanded to any other composition-resolved data sets, such as diurnal
853 changes of different compounds measured in ambient air, temporal changes of different
854 generations of species in a smog chamber, and composition-dependent size distributions. All of
855 the above data sets share a common property that the noise of the curve/spectrum is related to
856 the composition. Therefore, NSSC would facilitate the analysis by taking noise into consideration.](#)
857 Interpretation of the cluster-specific thermograms using frameworks such as that of
858 Schobesberger et al. (2018) will allow for more comprehensive understanding of the

859 thermochemical properties of the organic aerosol, the subject of future work. This will provide
860 insights into the thermal behavior of organic aerosol and the relative contributions of thermally
861 stable (e.g., monomer) versus thermally unstable (e.g., dimers or oligomers) compounds, the
862 volatility distribution of the thermally stable compounds, and the T-dependent rate coefficients
863 for oligomer dissociation and formation.

864 **6. Data Availability**

865 All data and the NSSC algorithm used in this publication are archived in the UC DASH data
866 repository (Cappa et al., 2019). The NSSC algorithm is also available at GitHub
867 (<https://github.com/chrisCappa/NSSC>), with the version used for this publication available as Li
868 and Cappa (2019).

869 **7. Author Contributions**

870 ZL developed the NSSC algorithm. ELD, SS, CJG, FDL-H, JL, JES, and ZL performed
871 measurements. ELD and SS performed detailed data processing. ZL and CDC analyzed data and
872 wrote the manuscript, with contributions from all co-authors.

873 **8. Acknowledgements**

874 This work was supported by the National Science Foundation under Grant No. ATM-
875 1151062. The experimental work described here was supported by the U.S. Department of
876 Energy ASR grants DE-SC0011791 and DE-SC0018221. E.L.D. was supported by the National
877 Science Foundation Graduate Research Fellowship (grant no. DGE-1256082) and S.S. was
878 supported by the Academy of Finland (grant nos. 272041 and 310682). The SOAFFEE campaign
879 was done at Pacific Northwest National Laboratory, supported by the U.S. Department of Energy
880 (DOE) Office of Science, Office of Biological and Environmental Research, as part of the
881 Atmospheric Systems Research (ASR) program. PNNL is operated for DOE by Battelle Memorial
882 Institute under contract DE-AC05-76RL01830.

9. References

- 884 Abdalmogith, S. S., and Harrison, R. M.: The use of trajectory cluster analysis to examine the long-
885 range transport of secondary inorganic aerosol in the UK, *Atmos Environ*, 39, 6686-6695,
886 <https://doi.org/10.1016/j.atmosenv.2005.07.059>, 2005.
- 887 Beddows, D. C. S., Dall'Osto, M., and Harrison, R. M.: Cluster Analysis of Rural, Urban, and
888 Curbside Atmospheric Particle Size Data, *Environ Sci Technol*, 43, 4694-4700,
889 <https://doi.org/10.1021/es803121t>, 2009.
- 890 Cape, J. N., Methven, J., and Hudson, L. E.: The use of trajectory cluster analysis to interpret trace
891 gas measurements at Mace Head, Ireland, *Atmos Environ*, 34, 3651-3663,
892 [https://doi.org/10.1016/S1352-2310\(00\)00098-4](https://doi.org/10.1016/S1352-2310(00)00098-4), 2000.
- 893 Cappa, C. D., Li, Z., D'Ambro, E. L., Schobesberger, S., Shilling, J. E., Lopez-Hilfiker, F., Liu, J., Gaston,
894 C. J., and Thornton, J. A.: Initial application of the noise-sorted scanning clustering algorithm to
895 the analysis of composition-dependent organic aerosol thermal desorption measurements, UC
896 Davis Dash, Dataset, <https://doi.org/10.25338/B87S43>, 2019
- 897 D'Ambro, E. L., Lee, B. H., Liu, J. M., Shilling, J. E., Gaston, C. J., Lopez-Hilfiker, F. D., Schobesberger,
898 S., Zaveri, R. A., Mohr, C., Lutz, A., Zhang, Z. F., Gold, A., Surratt, J. D., Rivera-Rios, J. C., Keutsch,
899 F. N., and Thornton, J. A.: Molecular composition and volatility of isoprene photochemical
900 oxidation secondary organic aerosol under low- and high-NO_x conditions, *Atmospheric Chemistry
901 and Physics*, 17, 159-174, <https://doi.org/10.5194/acp-17-159-2017>, 2017.
- 902 D'Ambro, E. L., Schobesberger, S., Zaveri, R. A., Shilling, J. E., Lee, B. H., Lopez-Hilfiker, F. D., Mohr,
903 C., and Thornton, J. A.: Isothermal Evaporation of alpha-Pinene Ozonolysis SOA: Volatility, Phase
904 State, and Oligomeric Composition, *Acs Earth Space Chem*, 2, 1058-1067,
905 <https://doi.org/10.1021/acsearthspacechem.8b00084>, 2018.
- 906 D'Ambro, E. L., Schobesberger, S., Gaston, C. J., Lopez-Hilfiker, F. D., Lee, B. H., Liu, J., Zelenyuk,
907 A., Bell, D., Cappa, C. D., Helgestad, T., Li, Z., Guenther, A., Wang, J., Wise, M., Caylor, R., Surratt,
908 J. D., Riedel, T., Hyttinen, N., Salo, V. T., Hasan, G., Kurtén, T., Shilling, J. E., and Thornton, J. A.:
909 Chamber-based insights into the factors controlling IEPOX SOA yield, composition, and volatility,
910 *Atmos. Chem. Phys. Discuss.*, 2019, 1-20, <https://doi.org/10.5194/acp-2019-271>, 2019.
- 911 Faxon, C., Hammes, J., Le Breton, M., Pathak, R. K., and Hallquist, M.: Characterization of organic
912 nitrate constituents of secondary organic aerosol (SOA) from nitrate-radical-initiated oxidation
913 of limonene using high-resolution chemical ionization mass spectrometry, *Atmospheric
914 Chemistry and Physics*, 18, 5467-5481, <https://doi.org/10.5194/acp-18-5467-2018>, 2018.
- 915 Gaston, C. J., Quinn, P. K., Bates, T. S., Gilman, J. B., Bon, D. M., Kuster, W. C., and Prather, K. A.:
916 The impact of shipping, agricultural, and urban emissions on single particle chemistry observed
917 aboard the R/V Atlantis during CalNex, *J Geophys Res-Atmos*, 118, 5003-5017,
918 <https://doi.org/10.1002/jgrd.50427>, 2013.
- 919 Gaston, C. J., Lopez-Hilfiker, F. D., Whybrew, L. E., Hadley, O., McNair, F., Gao, H. L., Jaffe, D. A.,
920 and Thornton, J. A.: Online molecular characterization of fine particulate matter in Port Angeles,
921 WA: Evidence for a major impact from residential wood smoke, *Atmos Environ*, 138, 99-107,
922 <https://doi.org/10.1016/j.atmosenv.2016.05.013>, 2016.
- 923 Giorio, C., Tapparo, A., Dall'Osto, M., Harrison, R. M., Beddows, D. C. S., Di Marco, C., and Nemitz,
924 E.: Comparison of three techniques for analysis of data from an Aerosol Time-of-Flight Mass

925 Spectrometer, *Atmos Environ*, 61, 316-326, <https://doi.org/10.1016/j.atmosenv.2012.07.054>,
926 2012.

927 Goldstein, A. H., and Galbally, I. E.: Known and unexplored organic constituents in the earth's
928 atmosphere, *Environ Sci Technol*, 41, 1514-1521, <https://doi.org/10.1021/es072476p>, 2007.

929 Gonzalez, T. F.: Clustering to Minimize the Maximum Intercluster Distance, *Theor Comput Sci*, 38,
930 293-306, [https://doi.org/10.1016/0304-3975\(85\)90224-5](https://doi.org/10.1016/0304-3975(85)90224-5), 1985.

931 Hamilton, J. F., Webb, P. J., Lewis, A. C., Hopkins, J. R., Smith, S., and Davy, P.: Partially oxidised
932 organic components in urban aerosol using GCXGC-TOF/MS, *Atmospheric Chemistry and Physics*,
933 4, 1279-1290, <https://doi.org/10.5194/acp-4-1279-2004>, 2004.

934 Huang, W., Saathoff, H., Pajunoja, A., Shen, X. L., Naumann, K. H., Wagner, R., Virtanen, A., Leisner,
935 T., and Mohr, C.: alpha-Pinene secondary organic aerosol at low temperature: chemical
936 composition and implications for particle viscosity, *Atmospheric Chemistry and Physics*, 18, 2883-
937 2898, <https://doi.org/10.5194/acp-18-2883-2018>, 2018.

938 Isaacman-VanWertz, G., Massoli, P., O'Brien, R. E., Nowak, J. B., Canagaratna, M. R., Jayne, J. T.,
939 Worsnop, D. R., Su, L., Knopf, D. A., Misztal, P. K., Arata, C., Goldstein, A. H., and Kroll, J. H.: Using
940 advanced mass spectrometry techniques to fully characterize atmospheric organic carbon:
941 current capabilities and remaining gaps, *Faraday Discussions*, 200, 579-598,
942 <https://doi.org/10.1039/c7fd00021a>, 2017.

943 Joo, T., Rivera-Rios, J. C., Takeuchi, M., Alvarado, M. J., and Ng, N. L.: Secondary Organic Aerosol
944 Formation from Reaction of 3-Methylfuran with Nitrate Radicals, *Acs Earth Space Chem*,
945 <https://doi.org/10.1021/acsearthspacechem.9b00068>, 2019.

946 Kirchner, U., Vogt, R., Natzeck, C., and Goschnick, J.: Single particle MS, SNMS, SIMS, XPS, and
947 FTIR spectroscopic analysis of soot particles during the AIDA campaign, *Journal of Aerosol Science*,
948 34, 1323-1346, [https://doi.org/10.1016/S0021-8502\(03\)00362-8](https://doi.org/10.1016/S0021-8502(03)00362-8), 2003.

949 Le Breton, M., Psichoudaki, M., Hallquist, M., Watne, A. K., Lutz, A., and Hallquist, A. M.:
950 Application of a FIGAERO ToF CIMS for on-line characterization of real-world fresh and aged
951 particle emissions from buses, *Aerosol Science and Technology*, 53, 244-259,
952 <https://doi.org/10.1080/02786826.2019.1566592>, 2019.

953 Lee, A. K. Y., Willis, M. D., Healy, R. M., Onasch, T. B., and Abbatt, J. P. D.: Mixing state of
954 carbonaceous aerosol in an urban environment: single particle characterization using the soot
955 particle aerosol mass spectrometer (SP-AMS), *Atmospheric Chemistry and Physics*, 15, 1823-
956 1841, <https://doi.org/10.5194/acp-15-1823-2015>, 2015.

957 Lee, B., Lopez-Hilfiker, F. D., D'Ambro, E. L., Zhou, P. T., Boy, M., Petaja, T., Hao, L. Q., Virtanen,
958 A., and Thornton, J. A.: Semi-volatile and highly oxygenated gaseous and particulate organic
959 compounds observed above a boreal forest canopy, *Atmospheric Chemistry and Physics*, 18,
960 11547-11562, <https://doi.org/10.5194/acp-18-11547-2018>, 2018.

961 Lee, B. H., Lopez-Hilfiker, F. D., Mohr, C., Kurten, T., Worsnop, D. R., and Thornton, J. A.: An Iodide-
962 Adduct High-Resolution Time-of-Flight Chemical-Ionization Mass Spectrometer: Application to
963 Atmospheric Inorganic and Organic Compounds, *Environ Sci Technol*, 48, 6309-6317,
964 <https://doi.org/10.1021/es500362a>, 2014.

965 Lee, B. H., Mohr, C., Lopez-Hilfiker, F. D., Lutz, A., Hallquist, M., Lee, L., Romer, P., Cohen, R. C.,
966 Iyer, S., Kurten, T., Hu, W. W., Day, D. A., Campuzano-Jost, P., Jimenez, J. L., Xu, L., Ng, N. L., Guo,
967 H. Y., Weber, R. J., Wild, R. J., Brown, S. S., Koss, A., de Gouw, J., Olson, K., Goldstein, A. H., Seco,
968 R., Kim, S., McAvey, K., Shepson, P. B., Starn, T., Baumann, K., Edgerton, E. S., Liu, J. M., Shilling,

969 J. E., Miller, D. O., Brune, W., Schobesberger, S., D'Ambro, E. L., and Thornton, J. A.: Highly
970 functionalized organic nitrates in the southeast United States: Contribution to secondary organic
971 aerosol and reactive nitrogen budgets, *P Natl Acad Sci USA*, 113, 1516-1521,
972 <https://doi.org/10.1073/pnas.1508108113>, 2016.

973 Li, Z., and Cappa, C. D.: Noise Sorted Scanning Clustering Algorithm (Version v1.0.3), Zenodo,
974 <https://doi.org/10.5281/zenodo.3361797>, 2019

975 Liu, J. M., D'Ambro, E. L., Lee, B. H., Lopez-Hilfiker, F. D., Zaveri, R. A., Rivera-Rios, J. C., Keutsch,
976 F. N., Iyer, S., Kurten, T., Zhang, Z. F., Gold, A., Surratt, J. D., Shilling, J. E., and Thornton, J. A.:
977 Efficient Isoprene Secondary Organic Aerosol Formation from a Non-IEPDX Pathway, *Environ Sci*
978 *Technol*, 50, 9872-9880, <https://doi.org/10.1021/acs.est.6b01872>, 2016.

979 Liu, S., Shilling, J. E., Song, C., Hiranuma, N., Zaveri, R. A., and Russell, L. M.: Hydrolysis of
980 Organonitrate Functional Groups in Aerosol Particles, *Aerosol Science and Technology*, 46, 1359-
981 1369, <https://doi.org/10.1080/02786826.2012.716175>, 2012.

982 Liu, S., Russell, L. M., Sueper, D. T., and Onasch, T. B.: Organic particle types by single-particle
983 measurements using a time-of-flight aerosol mass spectrometer coupled with a light scattering
984 module, *Atmospheric Measurement Techniques*, 6, 187-197, [https://doi.org/10.5194/amt-6-](https://doi.org/10.5194/amt-6-187-2013)
985 [187-2013](https://doi.org/10.5194/amt-6-187-2013), 2013.

986 Lopez-Hilfiker, F. D., Mohr, C., Ehn, M., Rubach, F., Kleist, E., Wildt, J., Mentel, T. F., Lutz, A.,
987 Hallquist, M., Worsnop, D., and Thornton, J. A.: A novel method for online analysis of gas and
988 particle composition: description and evaluation of a Filter Inlet for Gases and AEROSols
989 (FIGAERO), *Atmospheric Measurement Techniques*, 7, 983-1001, [https://doi.org/10.5194/amt-](https://doi.org/10.5194/amt-7-983-2014)
990 [7-983-2014](https://doi.org/10.5194/amt-7-983-2014), 2014.

991 Lopez-Hilfiker, F. D., Mohr, C., Ehn, M., Rubach, F., Kleist, E., Wildt, J., Mentel, T. F., Carrasquillo,
992 A. J., Daumit, K. E., Hunter, J. F., Kroll, J. H., Worsnop, D. R., and Thornton, J. A.: Phase partitioning
993 and volatility of secondary organic aerosol components formed from α -pinene ozonolysis and OH
994 oxidation: the importance of accretion products and other low volatility compounds,
995 *Atmospheric Chemistry and Physics*, 15, 7765-7776, [https://doi.org/10.5194/acp-15-7765-](https://doi.org/10.5194/acp-15-7765-2015)
996 [2015](https://doi.org/10.5194/acp-15-7765-2015), 2015.

997 Lopez-Hilfiker, F. D., Mohr, C., D'Ambro, E. L., Lutz, A., Riedel, T. P., Gaston, C. J., Iyer, S., Zhang,
998 Z., Gold, A., Surratt, J. D., Lee, B. H., Kurten, T., Hu, W. W., Jimenez, J., Hallquist, M., and Thornton,
999 J. A.: Molecular Composition and Volatility of Organic Aerosol in the Southeastern U.S.:
1000 Implications for IEPOX Derived SOA, *Environ Sci Technol*, 50, 2200-2209,
1001 <https://doi.org/10.1021/acs.est.5b04769>, 2016.

1002 Mohr, C., Lopez-Hilfiker, F. D., Yli-Juuti, T., Heitto, A., Lutz, A., Hallquist, M., D'Ambro, E. L.,
1003 Rissanen, M. P., Hao, L. Q., Schobesberger, S., Kulmala, M., Mauldin, R. L., Makkonen, U., Sipila,
1004 M., Petaja, T., and Thornton, J. A.: Ambient observations of dimers from terpene oxidation in the
1005 gas phase: Implications for new particle formation and growth, *Geophysical Research Letters*, 44,
1006 2958-2966, <https://doi.org/10.1002/2017gl072718>, 2017.

1007 Murphy, D. M., Middlebrook, A. M., and Warshawsky, M.: Cluster analysis of data from the
1008 Particle Analysis by Laser Mass Spectrometry (PALMS) instrument, *Aerosol Science and*
1009 *Technology*, 37, 382-391, <https://doi.org/10.1080/02786820300971>, 2003.

1010 Pinero-Garcia, F., Ferro-Garcia, M. A., Chham, E., Cobos-Diaz, M., and Gonzalez-Rodelas, P.: A
1011 cluster analysis of back trajectories to study the behaviour of radioactive aerosols in the south-

1012 east of Spain, *J Environ Radioactiv*, 147, 142-152, <https://doi.org/10.1016/j.jenvrad.2015.05.029>,
1013 2015.

1014 Prasse, E., Otkjaer, R. V., Crouse, J. D., Hethcox, J. C., Stoltz, B. M., Kjaergaard, H. G., and
1015 Wennberg, P. O.: Atmospheric autoxidation is increasingly important in urban and suburban
1016 North America, *P Natl Acad Sci USA*, 115, 64-69, <https://doi.org/10.1073/pnas.1715540115>, 2018.

1017 Rebotier, T. P., and Prather, K. A.: Aerosol time-of-flight mass spectrometry data analysis: A
1018 benchmark of clustering algorithms, *Anal Chim Acta*, 585, 38-54,
1019 <https://doi.org/10.1016/j.aca.2006.12.009>, 2007.

1020 Reitz, P., Zorn, S. R., Trimborn, S. H., and Trimborn, A. M.: A new, powerful technique to analyze
1021 single particle aerosol mass spectra using a combination of OPTICS and the fuzzy c-means
1022 algorithm, *Journal of Aerosol Science*, 98, 1-14, <https://doi.org/10.1016/j.jaerosci.2016.04.003>,
1023 2016.

1024 Roth, A., Schneider, J., Klimach, T., Mertes, S., van Pinxteren, D., Herrmann, H., and Borrmann, S.:
1025 Aerosol properties, source identification, and cloud processing in orographic clouds measured by
1026 single particle mass spectrometry on a central European mountain site during HCCT-2010,
1027 *Atmospheric Chemistry and Physics*, 16, 505-524, <https://doi.org/10.5194/acp-16-505-2016>,
1028 2016.

1029 Schobesberger, S., D'Ambro, E. L., Lopez-Hilfiker, F. D., Mohr, C., and Thornton, J. A.: A model
1030 framework to retrieve thermodynamic and kinetic properties of organic aerosol from
1031 composition-resolved thermal desorption measurements, *Atmospheric Chemistry and Physics*,
1032 18, 14757-14785, <https://doi.org/10.5194/acp-18-14757-2018>, 2018.

1033 Song, X. H., Hopke, P. K., Ferguson, D. P., and Prather, K. A.: Classification of single particles
1034 analyzed by ATOFMS using an artificial neural network, *ART-2A, Anal Chem*, 71, 860-865,
1035 <https://doi.org/10.1021/ac9809682>, 1999.

1036 Stolzenburg, D., Fischer, L., Vogel, A. L., Heinritzi, M., Schervish, M., Simon, M., Wagner, A. C.,
1037 Dada, L., Ahonen, L. R., Amorim, A., Baccarini, A., Bauer, P. S., Baumgartner, B., Bergen, A., Bianchi,
1038 F., Breitenlechner, M., Brilke, S., Mazon, S. B., Chen, D. X., Dias, A., Draper, D. C., Duplissy, J.,
1039 Haddad, I., Finkenzeller, H., Frege, C., Fuchs, C., Garmash, O., Gordon, H., He, X., Helm, J.,
1040 Hofbauer, V., Hoyle, C. R., Kim, C., Kirkby, J., Kontkanen, J., Kuerten, A., Lampilahti, J., Lawler, M.,
1041 Lehtipalo, K., Leiminger, M., Mai, H., Mathot, S., Mentler, B., Molteni, U., Nie, W., Nieminen, T.,
1042 Nowak, J. B., Ojdanic, A., Onnela, A., Passananti, M., Petaja, T., Quelever, L. L. J., Rissanen, M. P.,
1043 Sarnela, N., Schallhart, S., Tauber, C., Tome, A., Wagner, R., Wang, M., Weitz, L., Wimmer, D.,
1044 Xiao, M., Yan, C., Ye, P., Zha, Q., Baltensperger, U., Curtius, J., Dommen, J., Flagan, R. C., Kulmala,
1045 M., Smith, J. N., Worsnop, D. R., Hansel, A., Donahue, N. M., and Winkler, P. M.: Rapid growth of
1046 organic aerosol nanoparticles over a wide tropospheric temperature range, *P Natl Acad Sci USA*,
1047 115, 9122-9127, <https://doi.org/10.1073/pnas.1807604115>, 2018.

1048 Takahama, S., Gilardoni, S., Russell, L. M., and Kilcoyne, A. L. D.: Classification of multiple types
1049 of organic carbon composition in atmospheric particles by scanning transmission X-ray
1050 microscopy analysis, *Atmos Environ*, 41, 9435-9451,
1051 <https://doi.org/10.1016/j.atmosenv.2007.08.051>, 2007.

1052 Wang, D. S., and Ruiz, L. H.: Chlorine-initiated oxidation of n-alkanes under high-NO_x conditions:
1053 insights into secondary organic aerosol composition and volatility using a FIGAERO-CIMS,
1054 *Atmospheric Chemistry and Physics*, 18, 15535-15553, [https://doi.org/10.5194/acp-18-15535-](https://doi.org/10.5194/acp-18-15535-2018)
1055 [2018](https://doi.org/10.5194/acp-18-15535-2018), 2018.

1056 Wegner, T., Hussein, T., Hameri, K., Vesala, T., Kulmala, M., and Weber, S.: Properties of aerosol
1057 signature size distributions in the urban environment as derived by cluster analysis, Atmos
1058 Environ, 61, 350-360, <https://doi.org/10.1016/j.atmosenv.2012.07.048>, 2012.

1059 Zhao, W. X., Hopke, P. K., and Prather, K. A.: Comparison of two cluster analysis methods using
1060 single particle mass spectra, Atmos Environ, 42, 881-892,
1061 <https://doi.org/10.1016/j.atmosenv.2007.10.024>, 2008.

1062 Zhao, Y., Thornton, J. A., and Pye, H. O. T.: Quantitative constraints on autoxidation and dimer
1063 formation from direct probing of monoterpene-derived peroxy radical chemistry, P Natl Acad Sci
1064 USA, 115, 12142-12147, <https://doi.org/10.1073/pnas.1812147115>, 2018.

1065 Zhou, L. M., Hopke, P. K., and Venkatachari, P.: Cluster analysis of single particle mass spectra
1066 measured at Flushing, NY, Anal Chim Acta, 555, 47-56, <https://doi.org/10.1016/j.aca.2005.08.061>,
1067 2006.

1068

1069

1070 **10. Tables**1071 **Table 1.** Details of SOA formation and chamber conditions for all the example SOA systems.

Exp #	Precursor		Oxidant		Seeds		UV	T (°C)	RH (%)	NO ^{#§} (ppb)	M _p ^{#&} (µg/m ³)	FIGAERO Operation [§]
	Type	Conc. [#] (ppb)	Type	Conc. ^{##} (ppm)	Type	D _p ^{##} (nm)						
1*	α-pinene	10	OH (H ₂ O ₂)	1.0	AS ^{&}	50	On	25	50	-	5.1	Normal
2	Δ-3-carene	10	OH (H ₂ O ₂)	0.25	AS	50	On	25	50	-	5.2	Normal
3a										5	8.3	
3b	α-pinene	10	OH (H ₂ O ₂)	1.0	AS	50	On	25	50	10	9.2	Normal
3c										25	9.1	
4a												Normal
4b												1 h wait
4c	α-pinene	10	O ₃	0.1	PS ^{&&}	50	Off	25	80	-	4.0	3 h wait
4d												6 h wait
4e												24 h wait

* Experiment #1 is a case study used to test the performances of different clustering algorithms

Conc. of precursors are the concentrations expected in the chamber with the absence of any chemistry

For OH, conc. refers to concentration of H₂O₂ injected into the chamber; for O₃, conc. refers to steady-state concentration of O₃ in the chamber during SOA formation

Seed particles are size-selected in all the experiments

#§ NO concentration refers to the targeted NO concentration when NO is injected into the chamber. The actual steady-state concentration of NO is lower than targeted. "-" indicates that no external NO is added to the chamber

#& M_p is the estimated mass concentration of particles including SOA and seeds measured by SMPS when the chamber is at steady-state, except for experiment 4 where M_p is the mass concentration of SOA only

§ Normal operation mode means the desorption process starts immediately after collection period. X h wait means that particles are isothermally diluted for X hours before the desorption process is initiated

& AS = ammonium sulfate

&& PS = potassium sulfate

1072

1073

1074 **Table 2.** Comparison of different clustering algorithms

Clustering Algorithms	k-means	k-medoids	Mean-shift	DBSCAN	FPclustering	NSSC
Assign all the members?	Yes	Yes	No	No	Yes	No
Identify single-member clusters?	No	No	Yes	No	No	Yes
Robust solution?	No	No	No	Yes	No	Yes
Controlled distance from the center of clusters?	No	No	Yes	No	No	Yes
Influence of noise?	large	large	small	small	large	Small
Key preset parameters	N_c	N_c	ϵ, N_{min}	ϵ	Initial seed	ϵ, N_{min}
Software used in this study	Igor	R	Python	Igor	Igor	Igor

1075
1076

1077 **Table 3.** Parameters and thresholds used for the data processing and noise-sorted scanning clustering for
 1078 all the example experiments.

Expt #	SOA type	Pre-processing						Clustering				
		N_{total}	$N_{anomalous}$	$N_{filtered}$	$f_{m,filtered}$	ζ_{ref}	$f_{m,ref}$	ε	N_c	$N_{c,one}$	$f_{m,unclustered}$	$R_{interClst}$
1	α -pinene + OH	298	4	188	7.5	0.021	0.67	2.6	11	0	0.00	2.01
2	Δ -3-carene + OH	298	5	183	9.3	0.019	0.57	2.1	9	1	0.27	2.36
3a	α -pinene + OH + NO	Single	6	204	15.3	0.025	0.55	2.2	9	1	1.52	2.06
3b			6	204	17.5	-	-	-	9	1	1.72	-
3c			6	204	21.0	-	-	-	9	1	2.27	-
3a		Multi	2	208	15.5	0.025	0.55	2.2	9	1	1.52	2.06
3b			3	195	12.6	0.027	0.54	2.3	10	1	1.29	2.10
3c			6	200	12.8	0.028	0.43	2.5	12	1	1.21	1.96
4a	α -pinene + O ₃	Single	10	185	11.5	0.025	0.42	2.2	10	2	0.67	2.28
4b			10	185	14.0	-	-	-	10	2	0.79	-
4c			10	185	14.0	-	-	-	10	2	0.84	-
4d			10	185	13.8	-	-	-	10	2	0.83	-
4e			10	185	17.6	-	-	-	10	2	0.82	-
4a			Multi	1	191	11.4	0.025	0.41	2.2	11	2	1.04
4b	0	210		16.5	0.044	0.41	3.3	8	4	0.00	2.02	
4c	5	205		14.3	0.048	0.42	3.1	9	2	1.06	1.66	
4d	3	203		12.8	0.055	0.39	3.3	8	1	2.50	1.80	
4e	3	213		16.1	0.053	0.41	3.4	7	2	0.98	1.97	

N_{total} – Total number of ions characterized by CIMS

$N_{anomalous}$ – Number of anomalous ions

$N_{filtered}$ – Number of ions filtered out from the following clustering due to high levels of noises

$f_{m,filtered}$ – Mass fraction of the ions filtered out due to high levels of noises, expressed in %

ζ_{ref} – Noise threshold. Ions with noise levels above this threshold are excluded from clustering

$f_{m,ref}$ – The threshold of mass contribution (%) to identify an ion as significant

ε – distance criterion

N_c – Number of clusters determined with two or more members

$N_{c,one}$ – Number of clusters determined with only one member

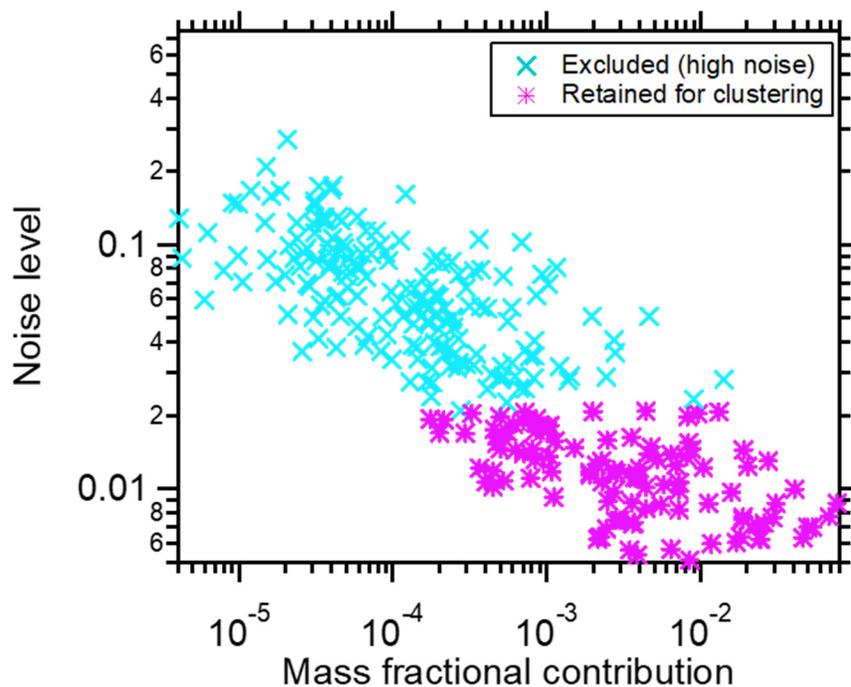
$f_{m,unclustered}$ – Mass fraction of unclustered ions, expressed in %

$R_{interClst}$ – The ratio of the average inter-cluster distance over the distance criterion ε

1079

1080

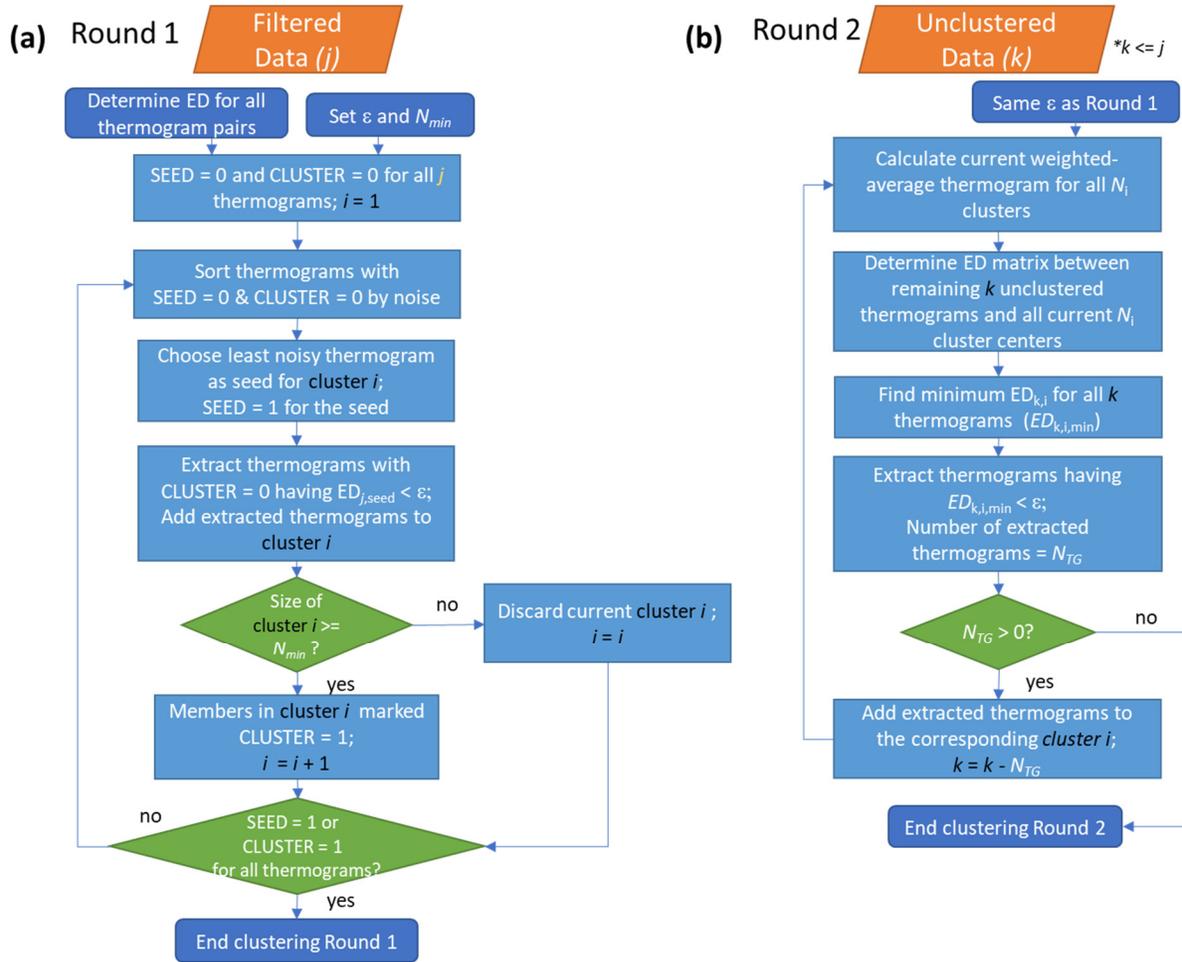
1081 **11. Figures**



1082

1083 **Figure 1:** The relationship between thermogram noise levels and the fractional contributions of the
1084 corresponding ions to total mass, for α -pinene + OH SOA. The noise threshold, $\zeta_{\text{ref}} = 0.021$ and is used to
1085 distinguish high-noise thermograms (cyan markers) from thermograms having acceptable noise levels
1086 (pink markers).

1087



1088

1089

1090

1091

1092

1093

1094

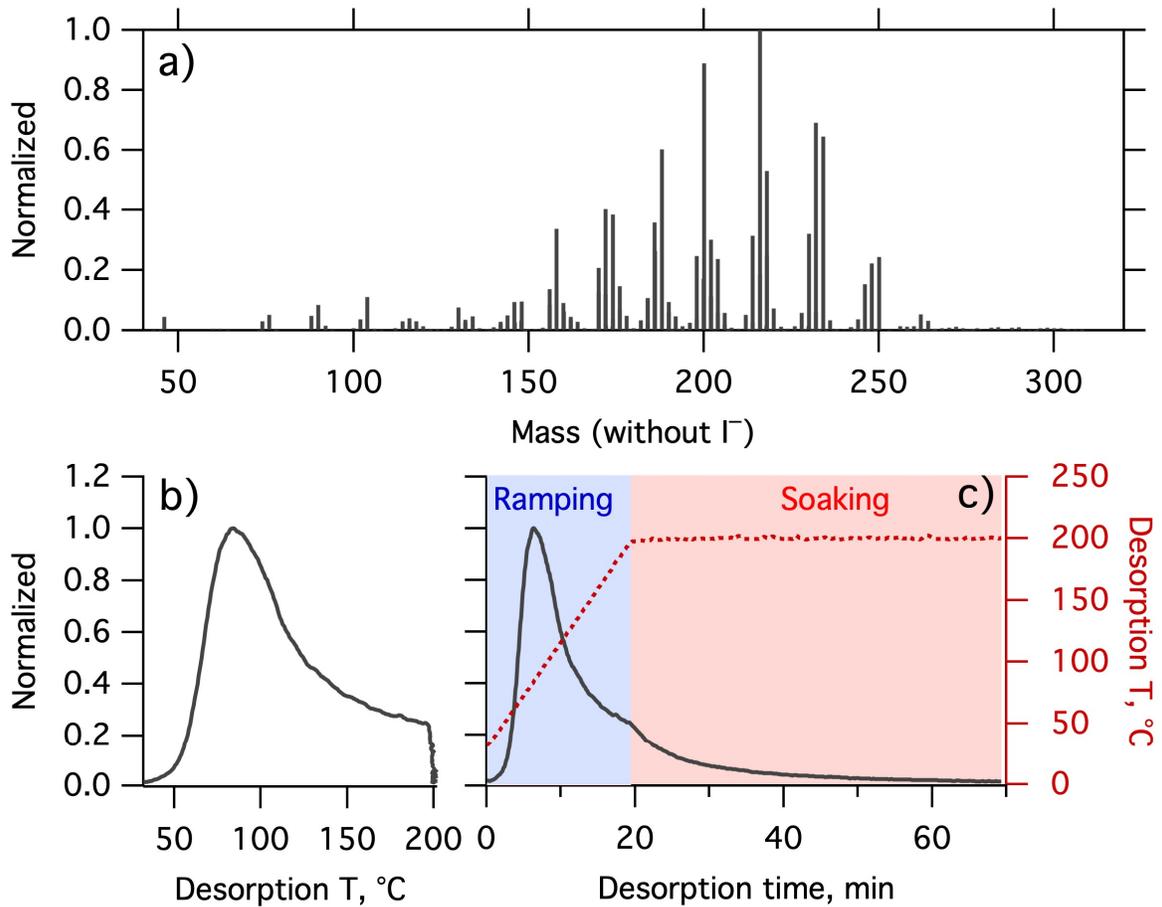
1095

1096

1097

1098

Figure 2: Flow of the noise-sorted scanning clustering. There are two rounds of clustering. (a) Round 1: The ED between all thermogram pairs are calculated and two parameters, ε and N_{min} , are set. Each thermogram is initialized with state SEED = 0 and CLUSTER = 0. Only thermograms with SEED = 0 and CLUSTER = 0 can serve as seeds, while thermograms with CLUSTER = 0 can be added to new clusters. The procedure terminates when all the thermograms are marked either SEED = 1 or CLUSTER = 1. (b) Round 2: Seeds are specified as the weighted-average thermogram for each cluster, and any remaining unclustered thermograms from Round 1 are potentially added to these clusters. With the indexing, j refers to the total number of thermograms, i to the number of clusters, and k to the number of unclustered thermograms after Round 1.



1100

1101

1102 **Figure 3.** (a) Mass spectrum of α -pinene + OH SOA measured by FIGAERO-CIMS. The mass excludes iodine.

1103 (b) Normalized thermogram of the bulk SOA versus temperature. (c) Normalized thermogram of the bulk

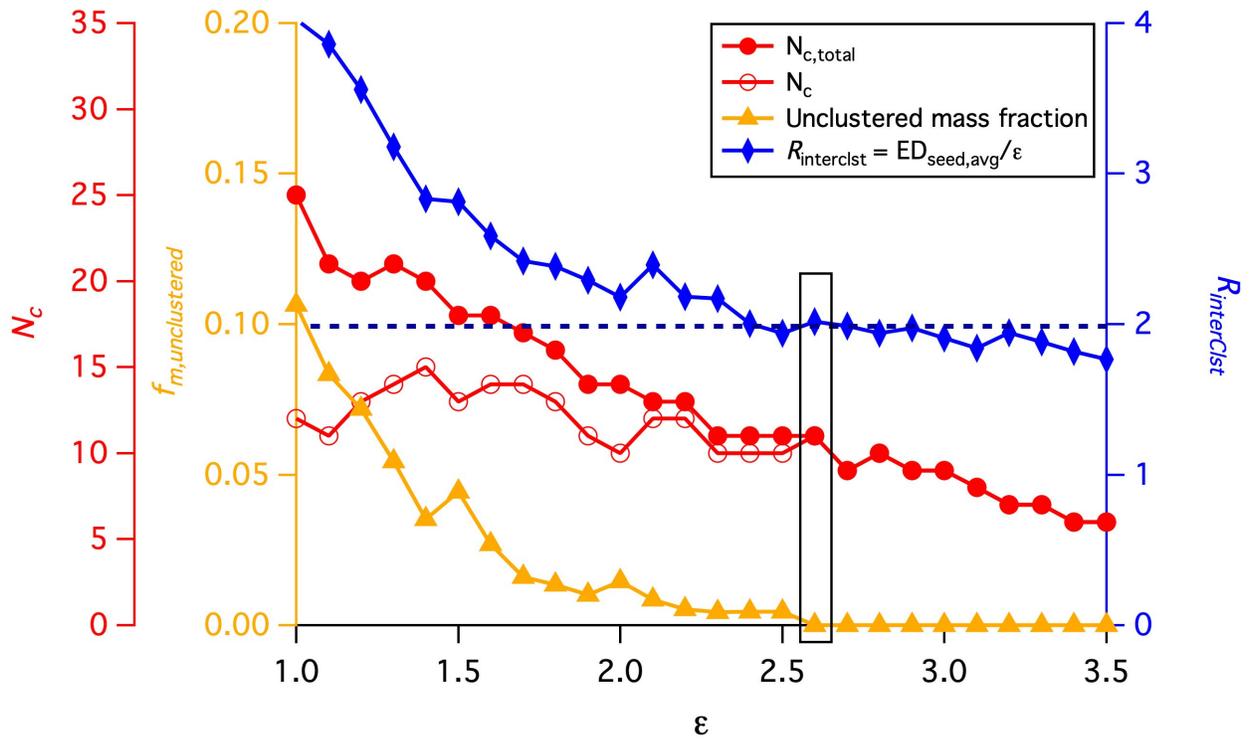
1104 SOA versus time (black line) and the variation in desorption temperature with time (dark red dashed line).

1105 The long tail during the soaking period is evident when the thermogram is considered in time space. The

1106 light blue shaded area denotes the ramping period and the pink shaded area the soaking period.

1107

1108



1109

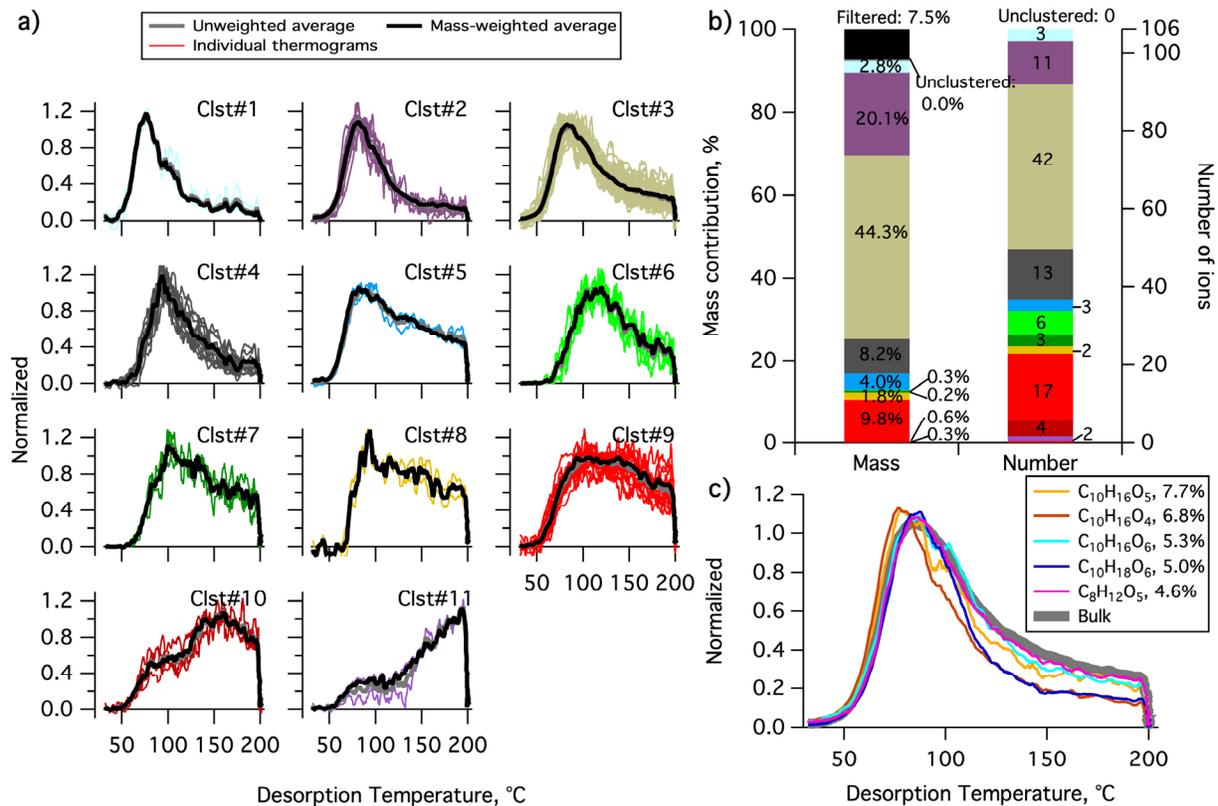
1110

Figure 4. The variation of four parameters, N_c , $N_{c,total}$, $f_{m,unclustered}$ and $R_{interClst}$ as a function of the distance criterion ϵ . The black horizontal dashed line guides the judgement for $R_{interClst} \geq 2$. The values highlighted by a rectangle are the values corresponding to the optimal ϵ used for the clustering analysis.

1111

1112

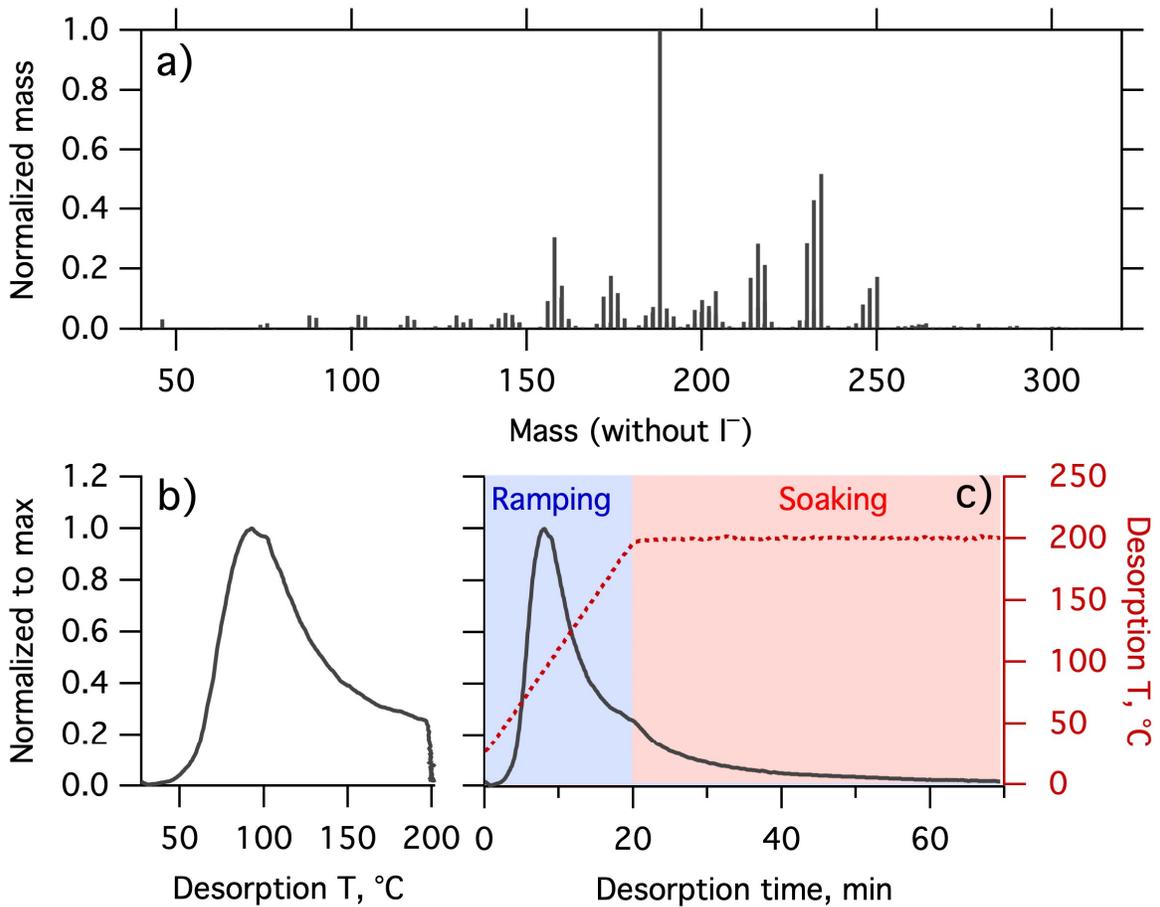
1113



1114
 1115 **Figure 5.** Clustering results for α -pinene + OH SOA. (a) Unweighted average thermograms (bold grey lines),
 1116 mass-weighted average thermograms (bold black lines) and individual members (colored lines) of the 11
 1117 clusters identified. (b) Percentage contribution of each cluster to the total mass, as well as the filtered out
 1118 and unclustered mass percentage (left bar), and the number of ions in each cluster and the unclustered
 1119 number of ions (right bar). (c) Thermograms of the top 5 ions in terms of mass contribution. The cluster
 1120 colors are consistent between (a) and (b).

1121

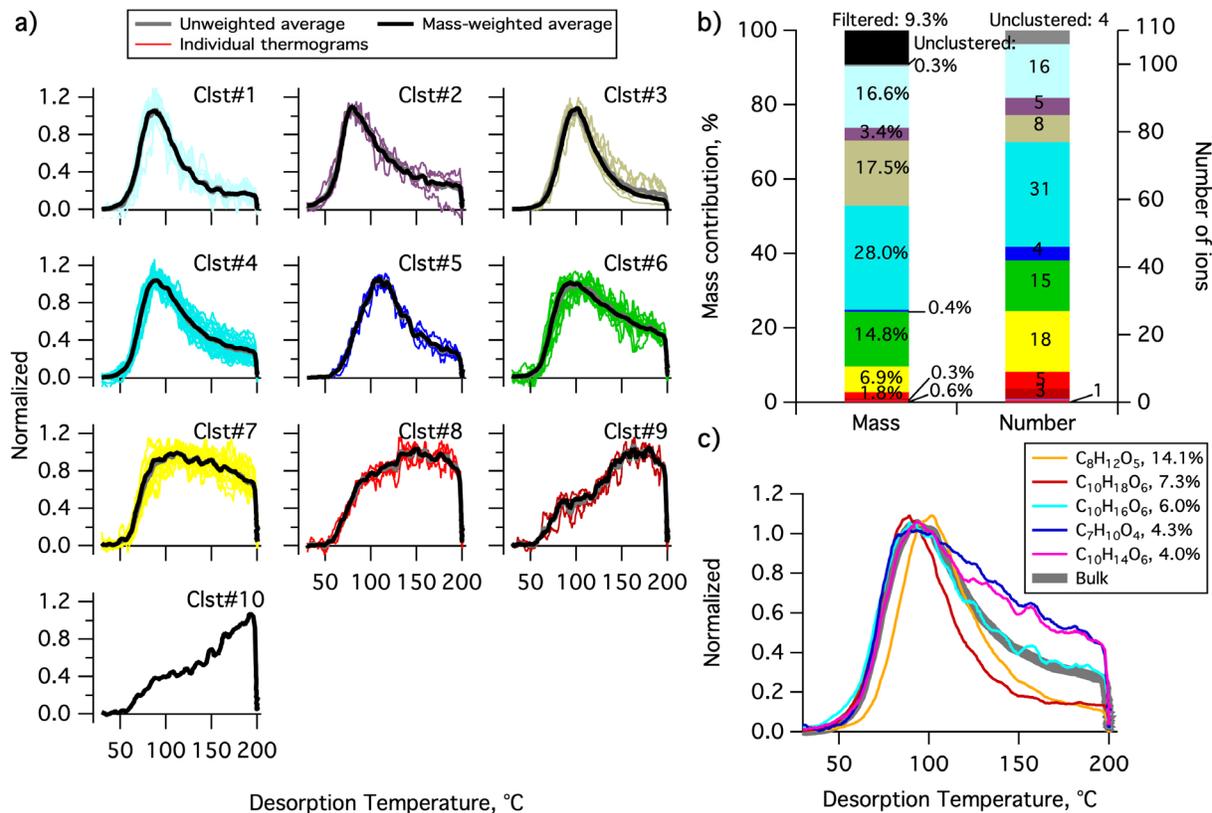
1122



1123

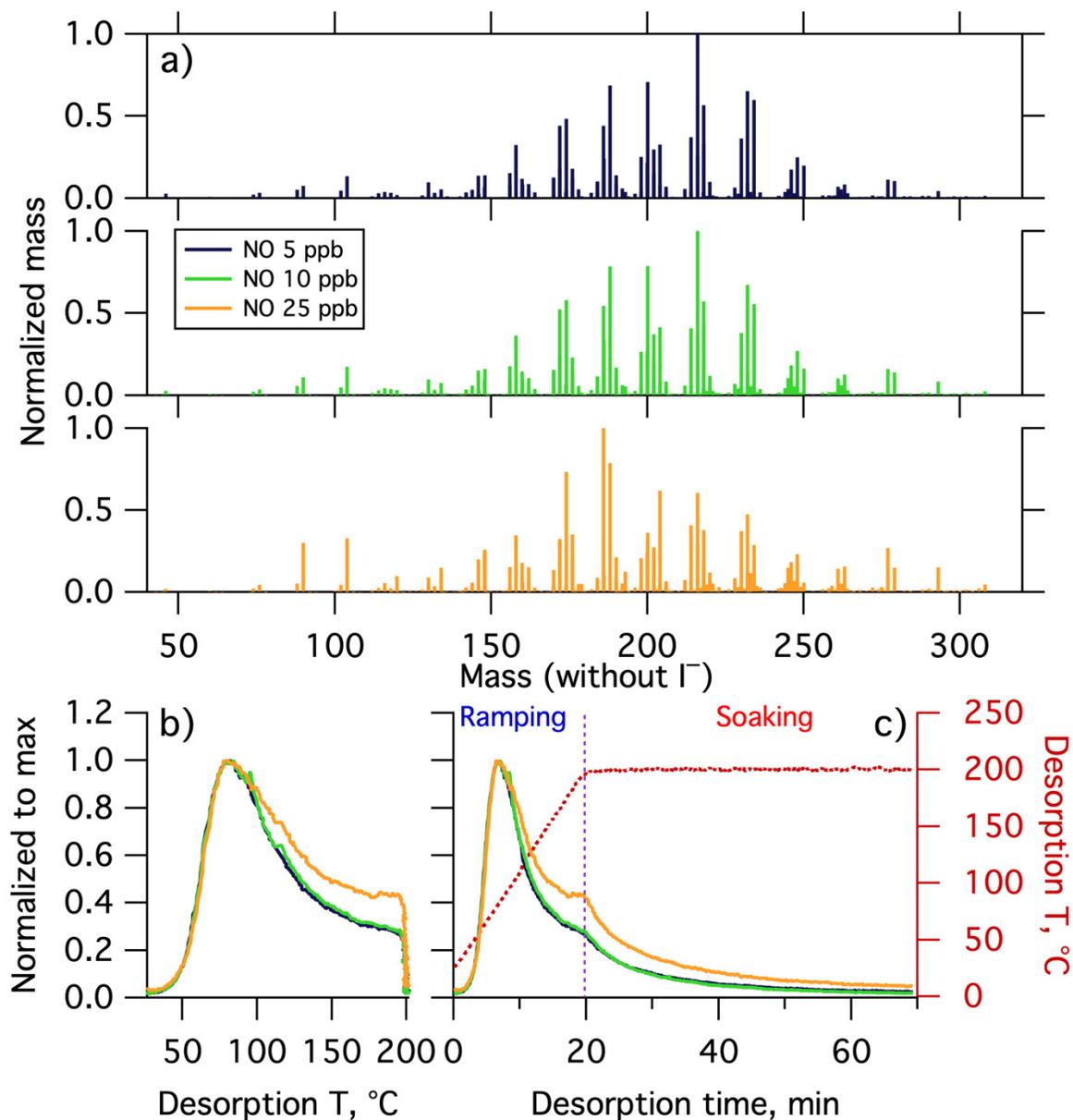
1124 **Figure 6.** Same as Figure 3, but for Δ -3-carene + OH SOA. (a) SOA mass spectrum measured by
1125 FIGAERO-CIMS. The mass excludes iodine. The normalized thermogram of the bulk SOA versus (b)
1126 temperature and (c) time. In (c) the light blue shaded area denotes the ramping period and the pink
1127 shaded area the soaking period.

1128



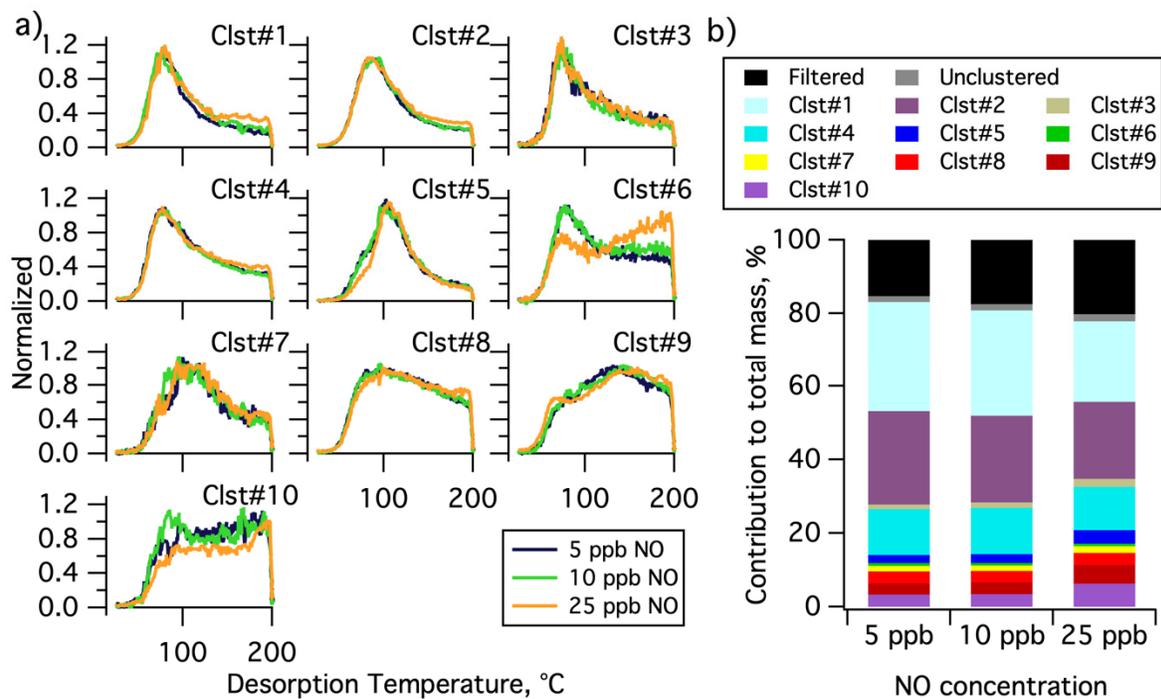
1129
 1130 **Figure 7.** Same as Figure 5, but for Δ -3-carene + OH SOA. (a) Unweighted average thermograms (bold grey
 1131 lines), mass-weighted average thermograms (bold black lines) and individual members (colored lines) of
 1132 the ten clusters identified. (b) Percentage contribution of each cluster to the total mass, as well as the
 1133 filtered out and unclustered mass percentage (left bar) and number of ions in each cluster and the
 1134 unclustered number of ions (right bar). (c) Thermograms of the top 5 ions in terms of mass contribution.
 1135 The cluster colors are consistent between (a) and (b).

1136



1137
 1138 **Figure 8.** (a) Mass spectra of α -pinene + OH SOA formed with different NO concentrations, normalized to
 1139 the most abundant ions mass concentration. The mass excludes iodine. Normalized thermograms of the
 1140 bulk SOA versus (b) temperature and (c) desorption time, with the desorption temperature shown in dark
 1141 red dashed line. The vertical purple dashed line delineates between ramping and soaking. In all the panels,
 1142 colors correspond to the NO concentration (see legend).

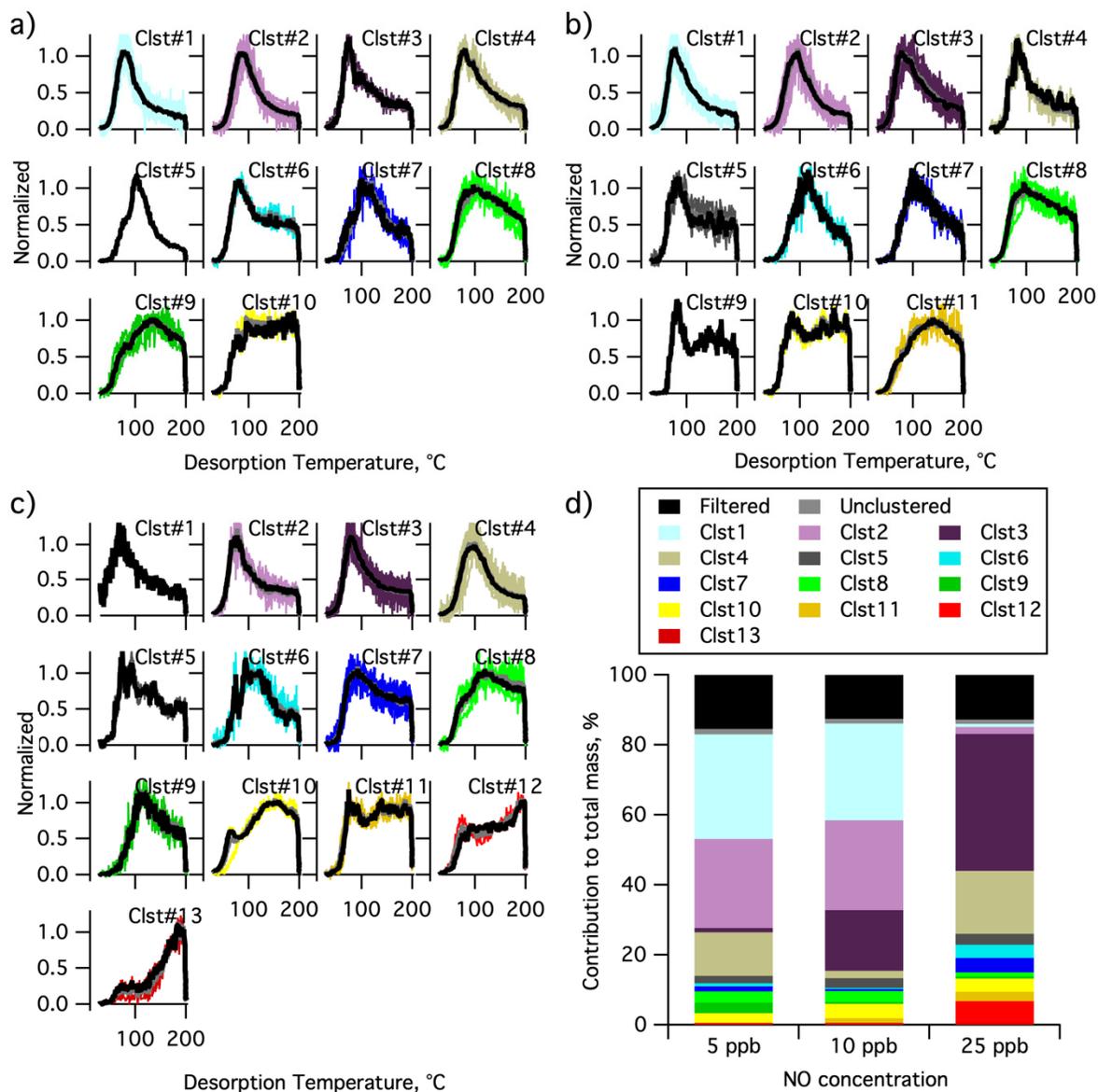
1143



1144

1145 **Figure 9.** Single clustering results for α -pinene + OH SOA as a function of NO concentration. (a)
 1146 Comparison of the normalized, weighted average thermograms of the ten clusters for the 5 ppb NO (navy),
 1147 10 ppb NO (green) and 25 ppb NO (orange) experiments. (b) Contribution of each cluster to the total mass,
 1148 including the contribution from filtered out ions (black) and unclustered ions (gray). The total mass is
 1149 calculated independently for each experiment.

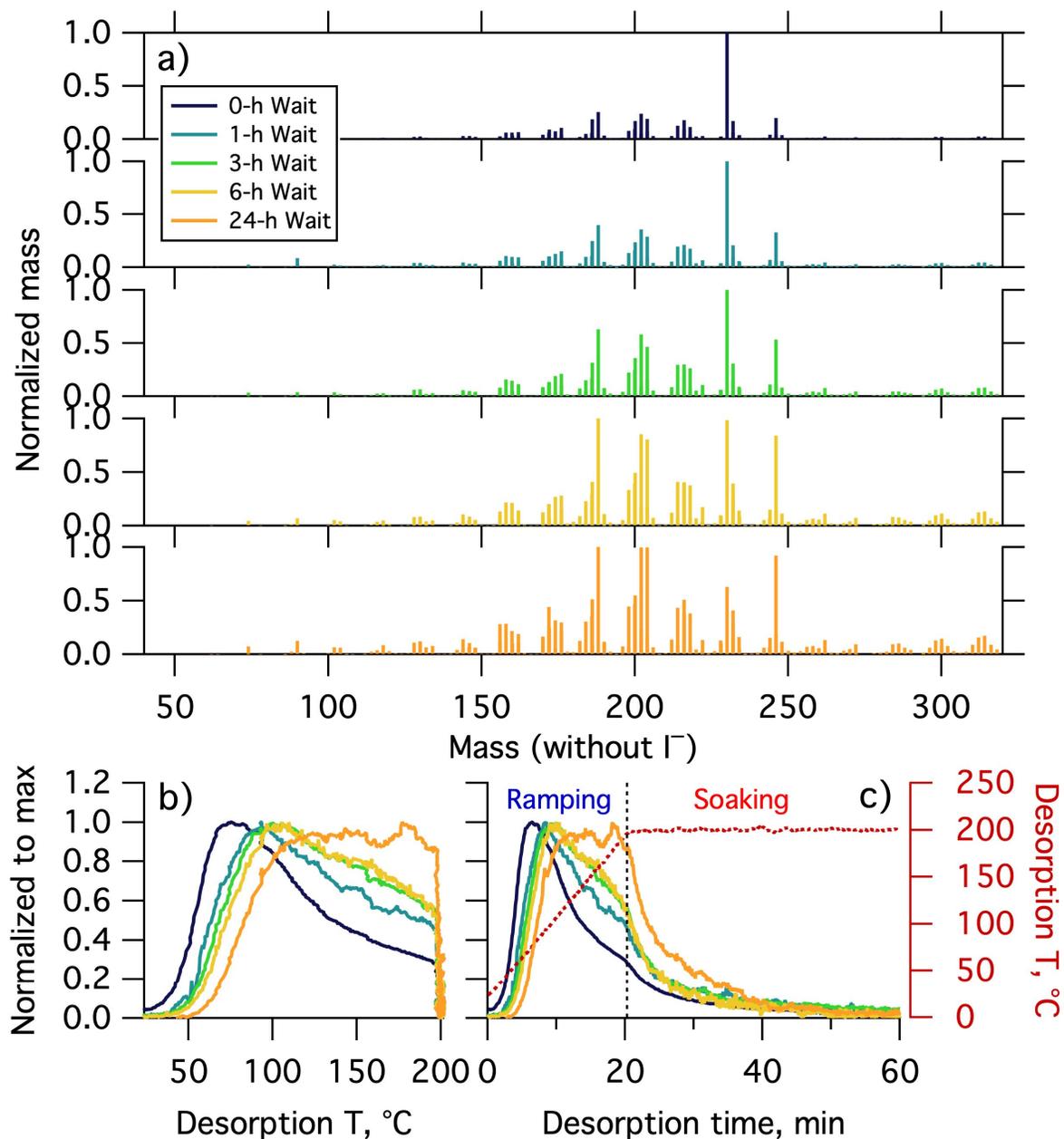
1150



1151

1152 **Figure 10.** Multiple clustering results for α -pinene + OH SOA as a function of NO concentration. Clustering
 1153 results are separately shown for the (a) 5 ppb NO, (b) 10 ppb NO, and (c) 25 ppb NO experiments. Each
 1154 panel includes unweighted average thermograms (grey lines), mass-weighted average thermograms
 1155 (black lines) and individual cluster members (colored lines). (d) Contribution of each cluster to the total
 1156 mass for each experiment. The mass contribution of filtered-out ions (black bar) and unclustered ions
 1157 (gray bar) are also shown.

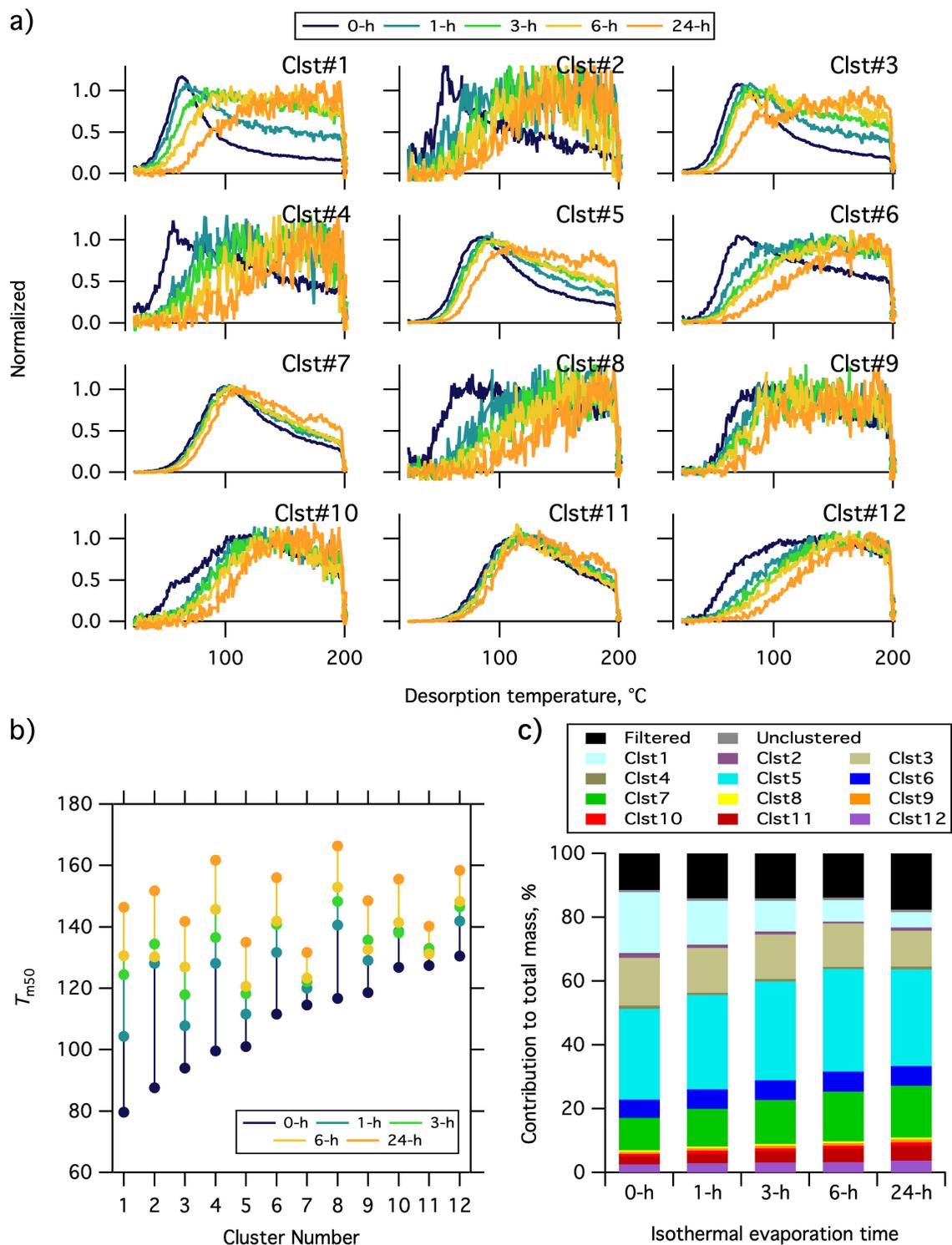
1158



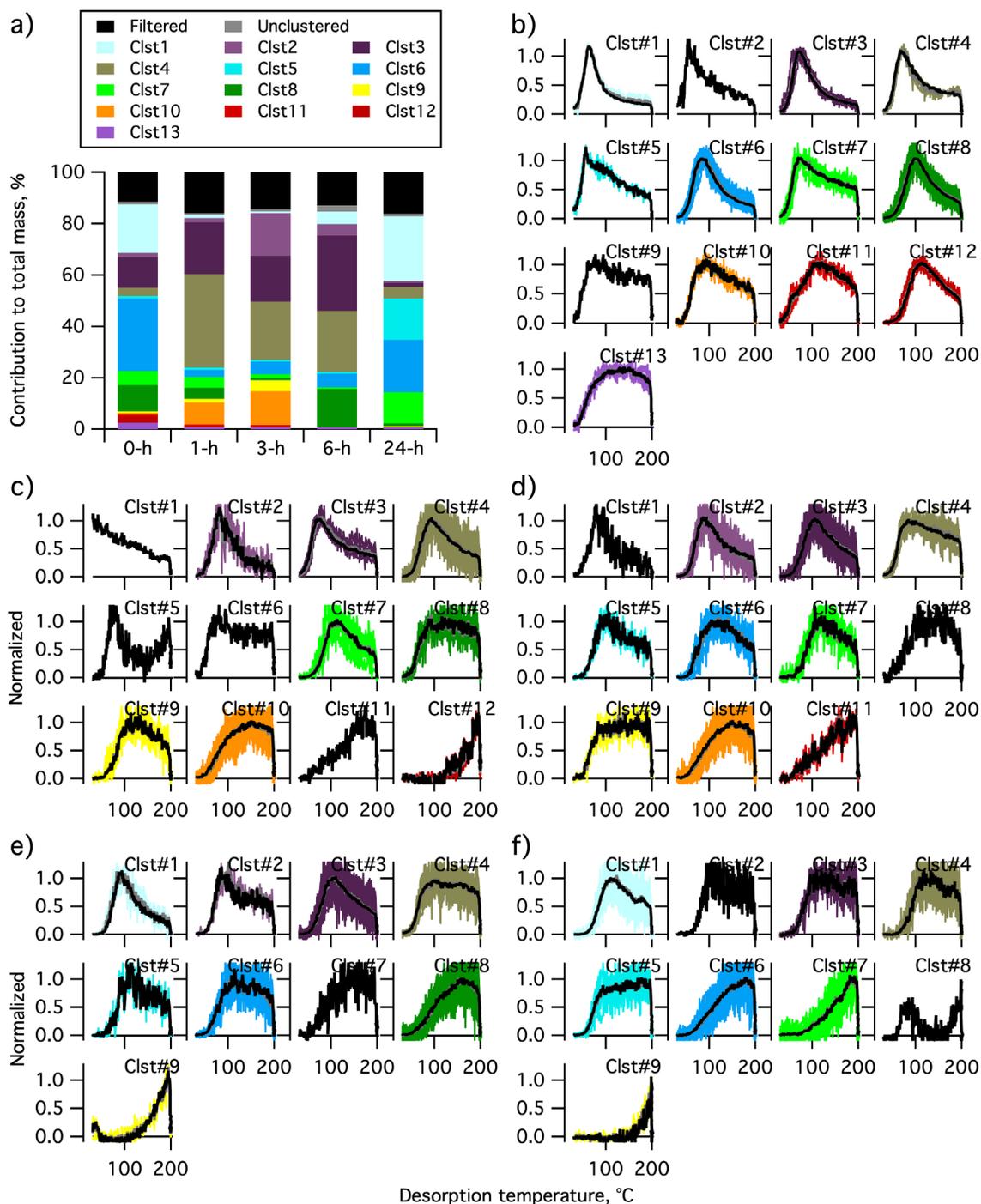
1160

1161 **Figure 11.** (a) Normalized mass spectra of α -pinene + O₃ SOA measured after different extents of
 1162 isothermal evaporation at room temperature. The mass excludes iodine. The normalized thermograms of
 1163 bulk SOA versus (b) temperature and (c) time, with the desorption temperature shown as a red dashed
 1164 line. The vertical black dashed line in (c) delineates between ramping and soaking. The mass spectrum or
 1165 thermogram colors indicate the isothermal evaporation time (see legend), with darker colors indicating
 1166 shorter times.

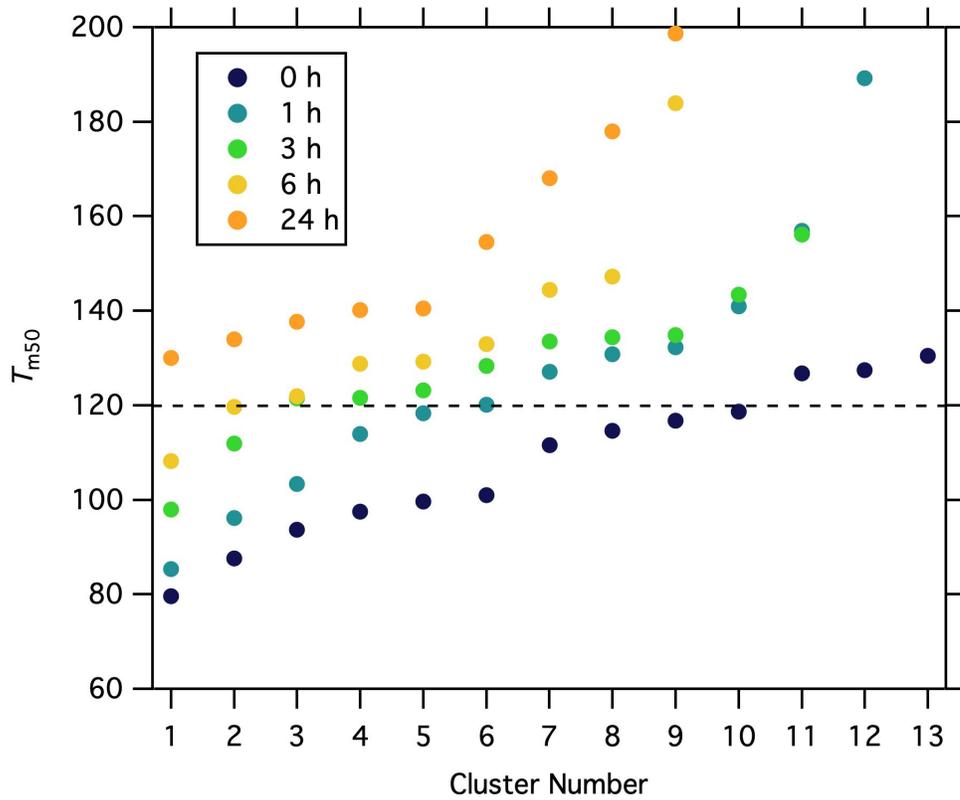
1167



1168
 1169 **Figure 12.** Single clustering results for α -pinene + O₃ SOA for different isothermal evaporation times. (a)
 1170 Comparison of the normalized, weighted-average thermograms of the 12 clusters of 0-h wait (navy), 1-h
 1171 wait (blue), 3-h wait (green), 6-h wait (yellow) and 24-h wait (orange) experiments. Note that the
 1172 absolute signals of all of the clusters decrease with evaporation, but to varying extents (Figure S6).



1173
 1174 **Figure 13.** Multiple clustering results for α -pinene + O₃ SOA as a function of isothermal evaporation time.
 1175 (a) Contribution of each cluster to the total mass for each experiment, along with the contributions of
 1176 filtered-out ions (black bar) and unclustered ions (gray bar). The number of clusters obtained generally
 1177 decreases with isothermal evaporation time. (b-f) The unweighted average (gray) and mass-weighted
 1178 average (black) thermograms, along with the thermograms of individual members of clusters for the (b)
 1179 0-h, (c) 1-h, (d) 3-h, (e) 6-h, and (f) 24-h wait experiments. The cluster colors are consistent between panels.



1180

1181 **Figure 14.** The T_{m50} values of the cluster-specific thermograms from multiple clustering for the five
 1182 isothermal evaporation experiments.

1183