

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

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

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

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

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

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

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

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

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

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

^c TofWerk AG, Thun, Switzerland

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

Abstract

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

40 of the thermograms of the mass spectral clusters will allow for more comprehensive
41 understanding of the thermochemical properties of OA particles.

42 **1. Introduction**

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

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

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

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

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

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

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

123 **2. Clustering Method Description**

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

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

128 **2.1. Data Preprocessing**

129 **2.1.1. Exclusion of anomalous thermograms**

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

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

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

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

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

164 **2.1.2. Euclidean Distance**

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

171

$$172 \quad ED_{a,b} = \sum_n \sqrt{(a_n - b_n)^2} \quad (1)$$

173

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

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

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

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

209 **2.1.3. Dealing with noise**

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

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

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

238 2.2. Noise-sorted Scanning Clustering (NSSC)

239 2.2.1. Algorithm description

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

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

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

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

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

285 **2.2.2. Post-clustering Processes**

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

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

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

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

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

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

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

328 2.2.3. Choosing the optimal ϵ

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

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

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

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

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

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

$$364 \quad R_{interClst} = \frac{ED_{seed,avg}}{\varepsilon} = \frac{\sum_{i=1}^{N_{c,total}} \sum_{j=1}^{N_{c,total}} ED_{seed,i,j}}{N_{c,total} \cdot (N_{c,total} - 1) \cdot \varepsilon} \quad (2)$$

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

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

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

382 **2.2.4. Summary**

383 The NSSC allows for clustering of ion peaks in temperature-dependent mass spectra
384 measured by the FIGAERO-CIMS, from which mass thermograms of the different clusters are
385 determined. The NSSC emphasizes contributions of ions having high signal-to-noise by selecting
386 seeds for the mass spectral clusters according to decreasing signal-to-noise. The NSSC also
387 accounts for the full temperature-dependent behavior of each ion, weighted towards the
388 temperature ramping period during which the ions generally exhibit more characteristic
389 desorption profiles. However, the NSSC requires as user input a distance criterion, ϵ , which
390 characterizes the minimum similarity required between a selected seed ion thermogram and all
391 other (non-clustered) ion-specific thermograms for the non-seed ion to be considered part of the
392 mass spectral cluster. The appropriate ϵ value must be uniquely determined for a given
393 experiment or set of experiments, but we recommend should be selected to provide both the
394 greatest amount of clustered mass and number of mass spectral clusters having two or more
395 members while also maintaining the greatest average separation between the mass spectral
396 clusters. Altogether, these steps facilitate robust, reproducible determination of mass spectral
397 clusters from a given data set.

398 **2.3. Alternative Clustering Methods**

399 We have alternatively considered the performance of some of the most commonly used
400 clustering algorithms (k-means, k-medoids, mean-shift, DBSCAN) and a less-commonly used one
401 (FPClustering (Gonzalez, 1985)) for interpreting FIGAERO-CIMS observations. The clustering
402 methods considered are summarized in **Table 2**, with some of their pros and cons listed, and
403 described in further detail in Appendix A. We discuss them briefly here in the context of FIGAERO-
404 CIMS data. All the methods considered require input of at least one key user-specified parameter.
405 These parameters and the associated clustering algorithms can be generally classified into two
406 categories: number-based and distance-based. Number-based clustering algorithms require

407 specifying the desired number of retrieved clusters; this includes k-means and k-medoids.
408 Number-based algorithms usually assign all members to clusters. The extent of similarity among
409 members of a cluster can vary greatly since there is no strict distance criterion for each cluster.
410 When applied to FIGAERO-CIMS thermograms, we have found these number-based algorithms
411 are particularly sensitive to the presence of noisy members and the initialization method. In
412 contrast, some clustering algorithms require specification of distance (similarity) criterion. This
413 includes the mean-shift, DBSCAN, and our NSSC algorithms. These distance-based algorithms
414 need not cluster all members of the initial population and generally emphasize intra-cluster
415 similarity or the density of the points. The methods differ in terms of the method used for
416 selection of the initial seed or center and the extent to which they emphasize point density versus
417 cluster similarity. Noisy members tend to naturally be excluded from any clusters. NSSC is a
418 variant of DBSCAN. It does, however, differ from the standard DBSCAN algorithm because NSSC
419 only searches for neighbors of the seed, while DBSCAN also searches for neighbors of the
420 neighbors. In doing so, NSSC emphasizes cluster similarity rather than point density. This is
421 particularly useful when clustering thermograms, as the behavior of the entire thermogram is
422 considered; inclusion of neighbors of neighbors may cluster together thermograms that exhibit
423 especially similar behavior in one region (e.g., the soaking period) but not another, an undesirable
424 result. Accordingly, the sorting of seeds by noise levels is a key aspect of the NSSC algorithm,
425 which we have found provides for more robust clustering results.

426 Most of these clustering algorithms, including k-means, k-medoids, and mean-shift, are
427 initialized with a random choice of the initial cluster centers (or seeds). For large data sets, this
428 randomness usually leads to different results of clustering with different runs. The extent to
429 which this impacts analysis and clustering of FIGAERO-CIMS data is considered using SOA from
430 the α -pinene + OH SOA system as the case study (Section 4.1). For the FIGAERO-CIMS data we
431 find that the various clustering results exhibit a moderate sensitivity to how the initial seeds are
432 selected for all of these algorithms, although the final clusters are generally similar between
433 different runs for the same input parameter. This may reflect either the relatively small size of
434 the data set (~300 members originally and ~100 members after noise screening) or that there are
435 generally characteristic peak shapes with overall good separation. However, some differences

436 between independent clustering runs result, which is undesirable. For FIGAERO-CIMS data we
437 know that not all thermograms are of equal quality, i.e. they have different noise levels reflecting
438 in part their different overall contributions to the total mass. The standard clustering methods
439 do not account for this information. The NSSC algorithm developed here takes into account this
440 measure of data quality and uses it to identify the seeds for clustering. This provides for an
441 entirely reproducible clustering and generally emphasizes the behavior of the ions that
442 contribute most to the FIGAERO-CIMS signal while still allowing for consideration of contributions
443 of low-signal ions.

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

453 In summary, we propose NSSC as the preferred algorithm in dealing with the FIGAERO data
454 set based on: (i) the ability to generate similar results as the other commonly used clustering
455 algorithms; (ii) good reproducibility and stability of results due to accounting for the noise of
456 individual thermograms; (iii) good control over the similarity within the clusters by using a
457 user-definable distance criterion; and (iv) a capability to identify unique thermograms as
458 one-member clusters.

459 **3. FIGAERO Measurements and Experiments**

460 **3.1. Instrument and experiment description**

461 The FIGAERO-CIMS instrument has been described previously in detail (Lee et al., 2014;
462 Lopez-Hilfiker et al., 2014). A brief description is provided here, with some additional details in
463 the Supplemental Material. The FIGAERO-CIMS measures the evolved gases from filter-collected

464 particles during temperature programmed thermal desorption. Thermal desorption of particles
465 occurs in two-stages: a “ramping” and “soaking” period. During ramping, the temperature
466 increases from room temperature to 200 °C, typically at 10 °C min⁻¹. Most OA mass desorbs
467 during the ramping stage. The temperature is held at 200 °C for ca. 30–40 mins during the soaking
468 period to facilitate evaporation of the remaining, low-volatility organic mass from the filter. The
469 evolved gas-phase compounds are measured using CIMS with the iodide (I⁻) reagent ion,
470 appropriate for characterization of generally highly oxygenated components comprising most
471 secondary organic aerosol (Lopez-Hilfiker et al., 2016; Isaacman-VanWertz et al., 2017; Lee et al.,
472 2018). The resulting signal or mass concentration versus temperature (or equivalently time)
473 curves for each ion constitute a thermogram. All individual thermograms are background
474 corrected by subtracting the observed thermograms from appropriate blank experiments. The
475 overall bulk thermogram is obtained by summing together the individual thermograms.

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

485 **3.2. General interpretation of FIGAERO-CIMS thermograms**

486 This work focuses on development of the clustering method, rather than on interpretation
487 of the FIGAERO-CIMS thermograms; an illustrative thermogram is shown in **Figure 3b**. However,
488 discussion of the clustering results is aided by a general understanding of how FIGAERO-CIMS
489 thermograms have been previously interpreted. Ions contributed by semi- and low-volatility
490 compounds that desorb directly tend to exhibit strongly peaked, Gaussian-like thermograms with
491 single-mode peaks between around 50 °C to 120 °C; the lower the peak desorption temperature
492 (T_{peak}) the higher the volatility of the desorbing compound (Lopez-Hilfiker et al., 2014; 2015). We

493 therefore refer to thermograms, or portions of thermograms, having this general shape as the
494 “monomeric” content of the ion hereafter; direct evaporation of thermally stable dimers or other
495 oligomers is possible, although will typically occur at higher temperatures due to the comparably
496 lower volatility of these compounds. When multiple monomeric compounds having different
497 vapor pressures contribute to the same ion, the resulting thermogram exhibits a broader peak
498 and shallower slopes or, in particular cases, multiple, distinct peaks (Lopez-Hilfiker et al., 2015).
499 However, very broad thermograms, especially those that peak at higher temperatures (> 120 °C
500 or so), can also indicate contributions from thermal decomposition of very low-volatility
501 monomers, dimers, and oligomers (Lopez-Hilfiker et al., 2015; Gaston et al., 2016; Schobesberger
502 et al., 2018). Dimers and oligomers can evaporate directly, without thermal decomposition, as
503 observed for isoprene-derived SOA (D'Ambro et al., 2017) and ambient monoterpene oxidation
504 products (Mohr et al., 2017). However, fragments of dimers or oligomers are generally more
505 abundant, indicating the importance of thermal decomposition for desorption of these low-
506 volatility compounds. Both direct evaporation of extremely low-volatility compounds and
507 decomposition of large molecules or oligomers can lead to high signal levels above ~120 °C. We
508 refer to both peaks and the slowly varying signal above ~120 °C as the “oligomeric” content of
509 the ion hereafter. We use the terms monomer and oligomer in a qualitative manner. A more
510 quantitative analysis of the thermograms can help distinguish between direct evaporation,
511 thermal decomposition, and the contributions of monomers versus oligomers (Schobesberger et
512 al., 2018), yet is beyond the scope of the current work.

513 **4. Example Applications**

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

519 4.1. α -pinene + OH SOA

520 A total of 298 ions were characterized by FIGAERO-CIMS for SOA generated from the
521 α -pinene + OH reaction (**Table 1**). Four ions were characterized as anomalous and excluded from
522 further analysis (see Section 2.1.1). The mass concentration of each ion was calculated by
523 integrating the signal across the entire desorption period and assuming an equal sensitivity of
524 CIMS for all the compounds. The total mass concentration is the sum of all the non-anomalous
525 ions. The mass spectrum and bulk thermogram of the remaining 294 ions are shown in **Figure 3**,
526 with the bulk thermogram shown versus both temperature (**Figure 3b**) and time (**Figure 3c**) to
527 illustrate the difference between the ramping and soaking periods. The individual thermograms
528 exhibited a variety of shapes. The noise threshold for this data set was $\xi_{\text{ref}} = 0.020893$. A total of
529 188 ions were screened out via noise filtering. The remaining 106 ions contribute 92.5% to the
530 total mass detected by FIGAERO-CIMS. The optimal ε was established through consideration of
531 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
532 set, we determine the optimal $\varepsilon = 2.6$. Choice of a much smaller ε , around 1.5, gives a maximum
533 in N_c , but leaves a large fraction of the mass unclustered. Choice of $\varepsilon = 2.1$ or 2.2 yields larger N_c
534 and $R_{\text{interClst}}$ than $\varepsilon = 2.6$, with a reasonably small $f_{m,\text{unclustered}}$. However, there is one type of
535 thermogram (Clst#11 in **Figure 5**) that is only captured with $\varepsilon \geq 2.6$ and this yields $f_{m,\text{unclustered}} = 0$.
536 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 $\varepsilon = 2.6$,
537 indicating increasing similarity between clusters with fewer types of shapes captured. The choice
538 of $\varepsilon = 2.6$ provides a compromise between maximizing N_c , minimizing $f_{m,\text{unclustered}}$, and keeping
539 $R_{\text{interClst}}$ above two. The parameters and thresholds used for this data set are summarized in **Table**
540 **3**.

541 A total of 11 clusters are identified with no one-member clusters. The unweighted and
542 mass-weighted average thermograms for each cluster are shown along with the thermograms of
543 individual members in **Figure 5a**. The differences between weighted and unweighted average
544 clusters are negligible, in general. Clusters are organized and numbered (as Clst# N) from low to
545 high T_{m50} , with deeper to shallower downslope. Clst#1 through Clst#6 all have a clear peak below
546 120 °C, but with different peak widths and downslopes. Clst#7 and Clst#8 are a bit noisier with

547 only a few members each, exhibiting a sharp upslope and shallow downslope. Clst#9 has a very
548 broad peak. Clst#10 peaks at around 150 °C after an initial rise and temporary plateau. Clst#11
549 exhibits behavior somewhat like Clst#10, but with a peak that occurs just into the soaking period,
550 evident if viewed in time space, at 200 °C with a rapid drop afterwards.

551 The total mass concentration of a given cluster ($M_{c,N}$) is the sum across all cluster members,
552 calculated by integrating the summed mass concentration across the entire desorption period.
553 The percentage mass contribution of each cluster, and of the unclustered and the noise-filtered
554 ions, as well as the number of members for each cluster are shown in **Figure 5b** and **Error!**
555 **Reference source not found.** Clst#2 and Clst#3 contain the majority of the mass (20.1% and
556 44.3%, respectively) and consist of nearly half of the clustered ions (11 and 42, respectively).
557 Clst#4 and Clst#9 also contain a notable percentage of the total mass (8.2% and 9.8%,
558 respectively) and include a notable number of ions (13 and 17, respectively). Other clusters
559 contribute relatively little to the total mass and contain a small fraction of ions.

560 The mass-weighted average molecular formulas ($C_xH_yO_zN_m$) differ between clusters, as do
561 the O:C and H:C atomic ratios (**Error! Reference source not found.**). There is no clear relationship
562 between T_{m50} (or cluster number) and the number of carbon atoms, MW, or O:C. There is,
563 however, a reasonable, inverse correlation between T_{m50} and H:C ($r^2 = 0.78$). The number of
564 carbon atoms is notably larger for Cluster 6 ($x = 11.1$) and Cluster 7 ($x = 15.3$); if those two clusters
565 are excluded there is an inverse relationship between T_{m50} and the number of carbon atoms (r^2
566 $= 0.79$) and with MW ($r^2 = 0.59$). While the reason for these two clusters having comparably large
567 numbers of carbon atoms is unknown, this nonetheless suggests that the contribution of
568 oligomer decomposition might increase for clusters having higher T_{m50} values.

569 Interpretation of previous FIGAERO-CIMS studies have largely focused on the behavior of
570 the bulk thermogram or of several major ions or sums of ions based on common factors such as
571 the number of carbon atoms (Lopez-Hilfiker et al., 2016; D'Ambro et al., 2017; D'Ambro et al.,
572 2018; Stolzenburg et al., 2018; Wang and Ruiz, 2018; Joo et al., 2019). The normalized
573 thermograms of the top five ions contributing most to the total mass for the experiments here
574 are shown in **Figure 5c**, along with the bulk thermogram. Together these five ions make up nearly
575 30% of the total mass, and exhibit very similar thermogram shapes to each other and to the bulk

576 thermogram and belong solely to either Clst#2 or Clst#3. Thus, examining these ions only would
577 capture only a fraction of the overall diversity in thermal behaviors. The clustering method
578 developed here provides a means to investigate more comprehensively the variability in volatility
579 between aerosol components.

580 **4.2. Δ -3-carene + OH SOA**

581 A total of 298 ions were characterized by FIGAERO-CIMS for SOA generated from the
582 reaction of Δ -3-carene + OH (**Table 1**). Five were identified as having anomalous thermograms
583 and excluded from further analysis. The mass spectrum and bulk thermograms of Δ -3-carene +
584 OH SOA are shown in **Figure 6**. Compared to the α -pinene +OH SOA described above, the mass
585 spectrum of Δ -3-carene SOA is quite different, with one ion ($C_8H_{12}O_5$) dominant. The bulk
586 thermograms of the two SOA systems both look bell-like, but with the Δ -3-carene SOA
587 thermogram having a peak temperature ca. 9 °C higher. After noise-filtering, 110 ions remained
588 for clustering, contributing 90.7% to the total mass. The optimal $\varepsilon = 2.1$, established again by
589 considering the system-specific dependence of N_c , $N_{c,one}$, $f_{m,unclustered}$ and $R_{interClst}$ on ε (Error!
590 Reference source not found.), with the parameters and thresholds summarized in **Table 3**.

591 Ten clusters are identified, including one one-member cluster, with thermograms shown in
592 **Figure 7a** and the mass contribution and number of ions in a cluster in **Figure 7b**. Chemical
593 properties of each cluster are summarized in Error! Reference source not found.. The general
594 characteristics of thermograms identified in the Δ -3-carene + OH SOA are similar to those of low-
595 NO_x α -pinene + OH SOA described above, but with different mass contributions. For example,
596 Clst#4 has nearly identical shape of the thermogram as Clst#3 in the α -pinene SOA but
597 contributes less to the total mass, 28.0% compared to 44.3%. Clst#6 in the Δ -3-carene SOA
598 contributes 14.8% to the total mass and resembles Clst#5 in the α -pinene SOA, which contributes
599 only 4.0% to the total mass.

600 In general, Clst#1 – 6 in the Δ -3-carene SOA all exhibit a peak below 120 °C, with clear peaks
601 of varying width and downslopes of varying steepness, but nominally in order of narrow to wide
602 and steep to shallow, respectively. These clusters carry the majority of the desorbed mass. Clst#7
603 and Clst#8 both exhibit relatively flat thermograms in the ramping period after their initial rise,

604 and contribute 9% to the total mass. Clst#9 has a peak temperature above 150 °C and Clst#10
605 reaches a maximum during the soaking period. These last two clusters contribute little to the
606 total mass (0.6% and 0.3%, respectively).

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

612 There are ultimately three major differences between the two SOA systems. For one, there
613 is a different relationship between fractional contribution and cluster number (and thus $T_{m,50}$)
614 between the two. Secondly, the α -pinene SOA contains ions with especially narrow peaks at ca.
615 100 °C (i.e., Clst#7 & 8), that are not observed with Δ -3-carene SOA (compare **Figure 5** with **Figure**
616 **7**). Lastly, the thermograms of the top five ions for Δ -3-carene SOA differ to a greater extent than
617 for α -pinene SOA. Although we are unable to determine the reasons for these differences here,
618 this illustrates the potential for clustering to help identify and understand differences between
619 different SOA systems.

620 **4.3. α -pinene + OH + NO SOA**

621 Thermograms from SOA generated from the reaction of α -pinene + OH at varying NO
622 concentrations (5 ppb, 10 ppb and 25 ppb; **Table 1**) are considered as a set of experiments.
623 Together, differences between them illustrate the impact of changes to the fate of RO₂ peroxy
624 radical intermediates on the SOA composition and thermal properties (Praske et al., 2018; Zhao
625 et al., 2018). Clustering proceeds here using two complementary approaches. In the single
626 clustering method, clustering is performed for one reference experiment (i.e., at one NO
627 concentration, 5 ppb, Expt#3a). Then, average thermograms are calculated for the other
628 experiments in the set using the same cluster members as identified in the reference experiment.
629 In the multiple clustering method, clusters are independently determined for each experiment in
630 the set, and the shapes, relative abundances, and contributing ions are compared between

631 experiments. For all three experiments, the same initial set of 298 ions were characterized by
632 FIGAERO-CIMS.

633 **4.3.1. Single Clustering**

634 The ions identified as anomalous in each experiment differed. This most likely results from
635 shifts in the background signal levels between experiments. To maintain consistency between
636 the three experiments, ions identified as anomalous in any of the experiments were excluded
637 from all the experiments, with four ions excluded in total. A total of 88 ions were kept for
638 clustering after noise-filtering using the 5 ppb NO reference experiment, contributing 84.5% to
639 the total mass. The optimal $\varepsilon = 2.2$ (Error! Reference source not found. and **Table 3**), resulting in
640 ten clusters with one one-member cluster. The same sets of ions were then used to calculate the
641 cluster-average thermograms for the 10 ppb and 25 ppb NO experiments. Chemical
642 characteristics of the clusters are summarized in Error! Reference source not found..

643 Mass spectra for the three experiments are compared in **Figure 8a** and the bulk
644 thermograms shown in **Figure 8b** and c. The 5 ppb NO and 10 ppb NO SOA mass spectra are
645 nearly identical. The mass spectrum for the 25 ppb NO experiment, however, exhibits a notable
646 shift of the most abundant ions towards lower m/z . The bulk thermograms for the 5 ppb and 10
647 ppb NO experiments are nearly identical, peaking near 80 °C. The 25 ppb NO bulk thermogram
648 similarly peaks near 80 °C, but exhibits a much slower decay as temperature increases further.
649 Additionally, the change in slope at the transition from the ramping to soaking period is more
650 pronounced in the 25 ppb NO experiment. Overall, a greater fraction of the mass desorbs above
651 100 °C and during the soaking period for the 25 ppb NO experiment compared to lower-NO
652 experiments.

653 Despite the differences in the bulk thermograms, the shapes of the weighted-average
654 thermograms of clusters for all the NO experiments are generally similar, with the exception of
655 Clst#6 (**Figure 9a**). In particular, the 25 ppb thermogram shape of Clst#6 differs substantially from
656 those of low-NO conditions, with a much reduced initial peak (around 80 °C) and an more
657 pronounced second peak at high temperature (around 200 °C). However, this cluster contributes
658 negligibly to the overall mass. There is some suggestion of similar behavior for Clst#10, although

659 to a lesser extent. For the three most abundant clusters, Clst#1, 2 and 4, there is a slightly
660 increased relative contribution of the 100-200 °C tail for 25 ppb NO, consistent with differences
661 in the bulk thermograms.

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

674 **4.3.2. Multiple Clustering**

675 With multiple clustering, each experiment was processed and clustered independently,
676 with experiment-specific ξ_{ref} , N_c , and ε , among other parameters (Error! Reference source not
677 found. and **Table 3**). The clustered thermograms from the three experiments are compared in
678 **Figure 10a-c**. The number of clusters identified increases with NO concentration. Comparison
679 between the shapes of the clusters from the 5 ppb NO (**Figure 10a**) and 10 ppb NO (**Figure 10b**)
680 experiments indicates generally similar types of thermograms, consistent with the single
681 clustering method. Ten of the 11 total 10 ppb clusters match with a 5 ppb cluster. The one
682 additional, unique cluster at 10 ppb NO (Clst#9), is a one-member cluster with a sharp, narrow
683 peak at low temperatures and a broader, shallow second peak at high temperatures. This ion was
684 filtered out due to high noise level in the 5 ppb NO experiment.

685 The 25 ppb NO experiment (**Figure 10c**) results in more clusters compared to the lower NO
686 experiments; 13 for the 25 ppb NO experiment versus 10 and 11 for the 5 and 10 ppb experiments,
687 respectively. Some of the 25 ppb NO clusters have shapes similar to the lower NO experiments,

688 but many differ substantially. For example, two of the unique 25 ppb NO clusters (Clst#12 and
689 #13) have thermograms for which the signal increases continuously through the ramping period
690 and even into the soaking period. These clusters were not found in the single clustering analysis
691 because the 5 ppb NO experiment was used as the reference.

692 The new types of thermograms observed in the 25 ppb NO experiment indicates either
693 formation of new compounds or a change in the relative contributions of different components
694 to the same ions. Either could result from a change in the fate of the peroxy radical intermediates
695 as the NO concentration increases, leading to notably different products. There were numerous
696 nitrogen-containing ions observed for the three experiments. These N-containing ions belong to
697 Clst#1 – 7 for all the three [NO] conditions (**Error! Reference source not found.**). The higher-
698 number clusters did not include N-containing ions, also indicating a limited influence of the
699 N-containing products on these lower-volatility thermograms, although fragmentation
700 complicates the interpretation. Overall, the formation of new N-containing compounds at the
701 high NO condition does not seem to explain the unique thermograms in the 25 ppb NO
702 experiments.

703 The percent contribution of different clusters to total mass, along with the noise-filtered
704 and unclustered ions, differ between experiments (**Figure 10d**). Note that for the multiple
705 clustering method, clusters having the same index number are not necessarily directly
706 comparable between experiments because different sets of ions are included. For example, while
707 Clst#1 in the 5 ppb and 10 ppb NO experiments are comparable, the most similar cluster in the
708 25 ppb experiment is Clst#2. Nonetheless, there are some common features shared by the same,
709 or closely indexed, clusters. For example, Clst#1 – 4 in all three experiments exhibit a narrow,
710 single peak with the peak temperature below 120 °C. The mass contribution of Clst#1 – 4 is similar
711 between the 5 and 10 ppb NO experiment, but ~15% lower in the 25 ppb NO experiment. Clusters
712 that reach their maximum signal at or above 150 °C (Clst#9, 10 for 5 ppb, Clst#10, 11 for 10 ppb
713 and Clst#10 – 13 for 25 ppb) together contribute ~6% in the low NO experiments and ~13% in
714 the high NO experiments. Thus, there is some evidence that at higher NO there is an increased
715 contribution of oligomeric compounds, indicated by the increased contribution of clusters that
716 peak at higher temperatures and exhibit broader overall thermograms. However, overall these

717 observations suggest complex shifts in the distribution of products, both monomeric and
718 oligomeric, with sufficient increases in NO to change the fate of the peroxy radical intermediates.

719 **4.4. α -pinene + O₃ SOA**

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

727 **4.4.1. Single Clustering**

728 Only a few ions, if any, were identified as anomalous in each experiment; a total of ten ions
729 were removed from all the experiments to maintain consistency between experiments. The mass
730 spectra and bulk thermograms of the remaining 302 ions for the five experiments are shown in
731 **Figure 11**. As the isothermal evaporation time increases, the mass spectrum changes significantly,
732 as previously reported by D'Ambro et al. (2018). In the no-wait experiment, the mass spectrum
733 is dominated by one ion, C₁₀H₁₄O₆. Upon isothermal evaporation, the relative abundance of this
734 ion notably decreases, with the extent of decrease increasing with wait time; over time, a greater
735 number of ions contribute to the total mass, both at lower and higher m/z . With isothermal
736 evaporation, the bulk thermograms also exhibit a shift from a more peaked shape, reminiscent
737 of that from a single compound (Lopez-Hilfiker et al., 2014), to a more flattened peak with a
738 shallower rise (**Figure 11**). In other words, with increasing isothermal evaporation the majority
739 of the mass desorbed during thermal desorption shifts from a lower to higher temperature region.
740 This behavior largely reflects the loss of comparably more volatile compounds during isothermal
741 evaporation, leaving behind SOA that is overall less volatile (**Error! Reference source not found.a**).
742 It can also in part be due to higher molecular weight, lower volatility compounds being produced
743 with time via accretion reactions in the condensed phase.

744 There are 12 clusters determined from the no-wait experiment, exhibiting a wide variety of
745 the shapes (**Figure 12a**), with the parameters used for data pre-processing and clustering
746 reported in **Table 3** and shown in Error! Reference source not found.. Focusing first on the no-
747 wait experiment, the cluster thermogram shapes include those having clear peaks at relatively
748 low temperatures (~ 60 °C) and with a sharp rise and fall (e.g., Clst#1-3), those having sharp peaks
749 at relatively low temperatures but with a shallow downward slope (e.g., Clst#6), those with a
750 broad peak at somewhat higher temperatures (~ 100 °C) and long tails (e.g., Clst#7), and those
751 having a wide peak at even higher temperatures ~ 120 °C with a very broad rise and fall (e.g.,
752 Clst#10).

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

768 Isothermal evaporation generally leads to a reduction of the monomeric character of
769 clusters, leaving behind components that exhibit increased oligomeric content. Differences in
770 how the individual cluster thermograms evolve with isothermal evaporation are therefore likely
771 indicative of differing relative contributions of monomeric versus oligomeric components. For
772 example, Clst#1 and Clst#10 have distinctly different shapes in the 0-h wait experiment, but very

773 similar shapes in the 24-h wait experiment. This indicates that ions in Clst#1 are not contributed
774 from a single component, as might be inferred from the single-mode peak in the 0-h wait
775 experiment. Instead, they are contributed by multiple components, though initially dominated
776 by monomeric compounds, so the shift in peak temperature and broadness is substantial. On the
777 other hand, ions in Clst#10 must also derive from multiple components, but with only a small
778 fraction of monomeric compounds that evaporate in the 24 hours. Consequently, the loss of
779 low-temperature mass is apparent yet small. In contrast, ions in clusters such as Clst#7 and 11
780 must be composed of only low-volatility components because they exhibit minimal changes in
781 the thermograms shapes.

782 The extent of mass loss with isothermal evaporation differs between clusters. In general,
783 clusters that exhibit larger changes in shape have greater total mass loss, although with variability
784 (**Error! Reference source not found.c**). Consequently, the mass contributions of the clusters
785 evolve with isothermal evaporation (**Figure 12b**). The contribution of Clst#1 decreases
786 significantly and most notably as wait time increases. The most prominent ion in the no-wait
787 experiment, $C_{10}H_{14}O_6$, is grouped in Clst#1. The continuous mass loss of Clst#1 indicates the rapid
788 evaporation of its members. The mass contributions of the other clusters that exhibited similar
789 changes in shape as Clst#1 (Clst#3, 5, and 6) remain comparably constant, although with Clst#3
790 decreasing slightly. The relative abundances of the clusters for which the thermograms shapes
791 changed negligibly (Clst#7 and 11) increase continually, implying of the slowest evaporation of
792 the ions in these two clusters in the 24-hr evaporation period.

793 For comparison, D'Ambro et al. (2018) reported changes in the shapes of the thermograms
794 for the five most abundant individual ions from the no-wait to 24-hr experiment, together
795 carrying ~15% of the particle mass. They observed the individual ion thermograms generally all
796 evolved in a manner similar to our Clst#1, 3 and 5, shifting from narrower, more peaked profiles
797 towards broader profiles with a shallower rise, less evident peak, and increased evaporation at
798 higher temperatures. Here, with the clustering of data, we are able to track the change of thermal
799 behaviors of ions carrying ~87% of the initial mass. We are able to confirm that ~70 % of the mass
800 exhibit similar thermal behaviors and responses to isothermal evaporation as the top five ions.
801 However, we are also able to identify another ~17% of the mass having initial thermograms not

802 characterized by the top five ions, including 12% of the mass (Clst#7 and 11) that behaves
803 distinctly different upon evaporation at room temperature.

804 **4.4.2. Multiple Clustering**

805 The number of clusters identified with the multiple-clustering method, using experiment-
806 specific optimal ε values (**Table 3** and Error! Reference source not found.), decreases with
807 isothermal evaporation time, from 13 (no-wait) to 12 (1 h) to 11 (3 h) and then to 9 (6 h and 24
808 h) (**Figure 13b-f**). The noise levels of the thermograms increase with evaporation time due to
809 decreasing absolute particle mass. Nonetheless, the typical shapes of the cluster-specific
810 thermograms clearly evolve with increasing isothermal evaporation. For short isothermal
811 evaporation times, many cluster-specific thermogram profiles are relatively narrow, peaking at
812 lower temperatures (70-120 °C) and with rapid rises and evident downslopes. For longer
813 isothermal evaporation times, the cluster-specific profiles instead have broad peaks with slow
814 rises and most of the mass desorbing at higher temperatures.

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

824 After isothermal evaporation, some cluster-specific thermograms have signals that increase
825 continuously during the ramping period, for example Clst#11 and 12 in the 1-h experiment; such
826 clusters were not observed in the no-wait experiment. The relative abundance of these very low-
827 volatility clusters increases with isothermal evaporation, from 1.7% in the 1-h experiment
828 (Clst#11 and 12) to 13.4% in the 24-hr experiment (Clst#7 and 9). The absence of these clusters
829 for the no-wait experiment suggests that they are formed over time through condensed-phase

830 reactions. Their increasing contribution over time may reflect both evaporation of higher
831 volatility components and continued formation. Clusters having thermograms with very broad
832 peaks, such as Clst#11 and 13 in the 0-h experiment are also observed in all the other experiments,
833 with increasing contribution to the total mass.

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

844 **5. Conclusions**

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

854 For the four different SOA systems we examined, more than 80% of the total mass is
855 clustered, with the number of clusters ranging from 9 to 13. The shapes of the cluster-specific
856 average thermograms exhibit substantial variation for a given system. Some have relatively sharp
857 peaks, others broad peaks with slowly decreasing signal as heating continues, and others still

858 having signals that continually increase up to very high temperatures or long desorption times.
859 The mass contribution of a cluster varies from 0.2% to 44.3%. A few (2-3) clusters usually contain
860 more than 50% of the total mass in all the chemical systems examined. Comparison of the cluster-
861 specific thermogram shapes between different SOA systems allows for qualitative assessment of
862 the similarity or uniqueness.

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

874 This paper focuses only on the description of the clustering algorithm and its potential as a
875 tool to characterize the thermal properties of organic aerosol in further detail. [The application of
876 NSSC can be potentially expanded to any other composition-resolved data sets, such as diurnal
877 changes of different compounds measured in ambient air, temporal changes of different
878 generations of species in a smog chamber, and composition-dependent size distributions. All of
879 the above data sets share a common property that the noise of the curve/spectrum is related to
880 the composition. Therefore, NSSC would facilitate the analysis by taking noise into consideration.](#)
881 Interpretation of the cluster-specific thermograms using frameworks such as that of
882 Schobesberger et al. (2018) will allow for more comprehensive understanding of the
883 thermochemical properties of the organic aerosol, the subject of future work. This will provide
884 insights into the thermal behavior of organic aerosol and the relative contributions of thermally
885 stable (e.g., monomer) versus thermally unstable (e.g., dimers or oligomers) compounds, the

886 volatility distribution of the thermally stable compounds, and the T-dependent rate coefficients
887 for oligomer dissociation and formation.

888 **6. Data Availability**

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

893 **7. Author Contributions**

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

897 **8. Acknowledgements**

898 This work was supported by the National Science Foundation under Grant No. ATM-
899 1151062. The experimental work described here was supported by the U.S. Department of
900 Energy ASR grants DE-SC0011791 and DE-SC0018221. E.L.D. was supported by the National
901 Science Foundation Graduate Research Fellowship (grant no. DGE-1256082) and S.S. was
902 supported by the Academy of Finland (grant nos. 272041 and 310682). The SOAFFEE campaign
903 was done at Pacific Northwest National Laboratory, supported by the U.S. Department of Energy
904 (DOE) Office of Science, Office of Biological and Environmental Research, as part of the
905 Atmospheric Systems Research (ASR) program. PNNL is operated for DOE by Battelle Memorial
906 Institute under contract DE-AC05-76RL01830.

907 **9. References**

908 Abdalmogith, S. S., and Harrison, R. M.: The use of trajectory cluster analysis to examine the long-
909 range transport of secondary inorganic aerosol in the UK, Atmos Environ, 39, 6686-6695,
910 <https://doi.org/10.1016/j.atmosenv.2005.07.059>, 2005.

911 Beddows, D. C. S., Dall'Osto, M., and Harrison, R. M.: Cluster Analysis of Rural, Urban, and
912 Curbside Atmospheric Particle Size Data, *Environ Sci Technol*, 43, 4694-4700,
913 <https://doi.org/10.1021/es803121t>, 2009.

914 Cape, J. N., Methven, J., and Hudson, L. E.: The use of trajectory cluster analysis to interpret trace
915 gas measurements at Mace Head, Ireland, *Atmos Environ*, 34, 3651-3663,
916 [https://doi.org/10.1016/S1352-2310\(00\)00098-4](https://doi.org/10.1016/S1352-2310(00)00098-4), 2000.

917 Cappa, C. D., Li, Z., D'Ambro, E. L., Schobesberger, S., Shilling, J. E., Lopez-Hilfiker, F., Liu, J., Gaston,
918 C. J., and Thornton, J. A.: Initial application of the noise-sorted scanning clustering algorithm to
919 the analysis of composition-dependent organic aerosol thermal desorption measurements, UC
920 Davis Dash, Dataset, <https://doi.org/10.25338/B87S43>, 2019

921 D'Ambro, E. L., Lee, B. H., Liu, J. M., Shilling, J. E., Gaston, C. J., Lopez-Hilfiker, F. D., Schobesberger,
922 S., Zaveri, R. A., Mohr, C., Lutz, A., Zhang, Z. F., Gold, A., Surratt, J. D., Rivera-Rios, J. C., Keutsch,
923 F. N., and Thornton, J. A.: Molecular composition and volatility of isoprene photochemical
924 oxidation secondary organic aerosol under low- and high-NO_x conditions, *Atmospheric Chemistry
925 and Physics*, 17, 159-174, <https://doi.org/10.5194/acp-17-159-2017>, 2017.

926 D'Ambro, E. L., Schobesberger, S., Zaveri, R. A., Shilling, J. E., Lee, B. H., Lopez-Hilfiker, F. D., Mohr,
927 C., and Thornton, J. A.: Isothermal Evaporation of alpha-Pinene Ozonolysis SOA: Volatility, Phase
928 State, and Oligomeric Composition, *Acs Earth Space Chem*, 2, 1058-1067,
929 <https://doi.org/10.1021/acsearthspacechem.8b00084>, 2018.

930 D'Ambro, E. L., Schobesberger, S., Gaston, C. J., Lopez-Hilfiker, F. D., Lee, B. H., Liu, J., Zelenyuk,
931 A., Bell, D., Cappa, C. D., Helgestad, T., Li, Z., Guenther, A., Wang, J., Wise, M., Caylor, R., Surratt,
932 J. D., Riedel, T., Hyttinen, N., Salo, V. T., Hasan, G., Kurtén, T., Shilling, J. E., and Thornton, J. A.:
933 Chamber-based insights into the factors controlling IEPOX SOA yield, composition, and volatility,
934 *Atmos. Chem. Phys. Discuss.*, 2019, 1-20, <https://doi.org/10.5194/acp-2019-271>, 2019.

935 Faxon, C., Hammes, J., Le Breton, M., Pathak, R. K., and Hallquist, M.: Characterization of organic
936 nitrate constituents of secondary organic aerosol (SOA) from nitrate-radical-initiated oxidation
937 of limonene using high-resolution chemical ionization mass spectrometry, *Atmospheric
938 Chemistry and Physics*, 18, 5467-5481, <https://doi.org/10.5194/acp-18-5467-2018>, 2018.

939 Gaston, C. J., Quinn, P. K., Bates, T. S., Gilman, J. B., Bon, D. M., Kuster, W. C., and Prather, K. A.:
940 The impact of shipping, agricultural, and urban emissions on single particle chemistry observed
941 aboard the R/V Atlantis during CalNex, *J Geophys Res-Atmos*, 118, 5003-5017,
942 <https://doi.org/10.1002/jgrd.50427>, 2013.

943 Gaston, C. J., Lopez-Hilfiker, F. D., Whybrew, L. E., Hadley, O., McNair, F., Gao, H. L., Jaffe, D. A.,
944 and Thornton, J. A.: Online molecular characterization of fine particulate matter in Port Angeles,
945 WA: Evidence for a major impact from residential wood smoke, *Atmos Environ*, 138, 99-107,
946 <https://doi.org/10.1016/j.atmosenv.2016.05.013>, 2016.

947 Giorio, C., Tapparo, A., Dall'Osto, M., Harrison, R. M., Beddows, D. C. S., Di Marco, C., and Nemitz,
948 E.: Comparison of three techniques for analysis of data from an Aerosol Time-of-Flight Mass
949 Spectrometer, *Atmos Environ*, 61, 316-326, <https://doi.org/10.1016/j.atmosenv.2012.07.054>,
950 2012.

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

953 Gonzalez, T. F.: Clustering to Minimize the Maximum Intercluster Distance, *Theor Comput Sci*, 38,
954 293-306, [https://doi.org/10.1016/0304-3975\(85\)90224-5](https://doi.org/10.1016/0304-3975(85)90224-5), 1985.

955 Hamilton, J. F., Webb, P. J., Lewis, A. C., Hopkins, J. R., Smith, S., and Davy, P.: Partially oxidised
956 organic components in urban aerosol using GCXGC-TOF/MS, *Atmospheric Chemistry and Physics*,
957 4, 1279-1290, <https://doi.org/10.5194/acp-4-1279-2004>, 2004.

958 Huang, W., Saathoff, H., Pajunoja, A., Shen, X. L., Naumann, K. H., Wagner, R., Virtanen, A., Leisner,
959 T., and Mohr, C.: alpha-Pinene secondary organic aerosol at low temperature: chemical
960 composition and implications for particle viscosity, *Atmospheric Chemistry and Physics*, 18, 2883-
961 2898, <https://doi.org/10.5194/acp-18-2883-2018>, 2018.

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

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

970 Kirchner, U., Vogt, R., Natzeck, C., and Goschnick, J.: Single particle MS, SNMS, SIMS, XPS, and
971 FTIR spectroscopic analysis of soot particles during the AIDA campaign, *Journal of Aerosol Science*,
972 34, 1323-1346, [https://doi.org/10.1016/S0021-8502\(03\)00362-8](https://doi.org/10.1016/S0021-8502(03)00362-8), 2003.

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

977 Lee, A. K. Y., Willis, M. D., Healy, R. M., Onasch, T. B., and Abbatt, J. P. D.: Mixing state of
978 carbonaceous aerosol in an urban environment: single particle characterization using the soot
979 particle aerosol mass spectrometer (SP-AMS), *Atmospheric Chemistry and Physics*, 15, 1823-
980 1841, <https://doi.org/10.5194/acp-15-1823-2015>, 2015.

981 Lee, B., Lopez-Hilfiker, F. D., D'Ambro, E. L., Zhou, P. T., Boy, M., Petaja, T., Hao, L. Q., Virtanen,
982 A., and Thornton, J. A.: Semi-volatile and highly oxygenated gaseous and particulate organic
983 compounds observed above a boreal forest canopy, *Atmospheric Chemistry and Physics*, 18,
984 11547-11562, <https://doi.org/10.5194/acp-18-11547-2018>, 2018.

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

989 Lee, B. H., Mohr, C., Lopez-Hilfiker, F. D., Lutz, A., Hallquist, M., Lee, L., Romer, P., Cohen, R. C.,
990 Iyer, S., Kurten, T., Hu, W. W., Day, D. A., Campuzano-Jost, P., Jimenez, J. L., Xu, L., Ng, N. L., Guo,
991 H. Y., Weber, R. J., Wild, R. J., Brown, S. S., Koss, A., de Gouw, J., Olson, K., Goldstein, A. H., Seco,
992 R., Kim, S., McAvey, K., Shepson, P. B., Starn, T., Baumann, K., Edgerton, E. S., Liu, J. M., Shilling,
993 J. E., Miller, D. O., Brune, W., Schobesberger, S., D'Ambro, E. L., and Thornton, J. A.: Highly
994 functionalized organic nitrates in the southeast United States: Contribution to secondary organic
995 aerosol and reactive nitrogen budgets, *P Natl Acad Sci USA*, 113, 1516-1521,
996 <https://doi.org/10.1073/pnas.1508108113>, 2016.

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

999 Liu, J. M., D'Ambro, E. L., Lee, B. H., Lopez-Hilfiker, F. D., Zaveri, R. A., Rivera-Rios, J. C., Keutsch,
1000 F. N., Iyer, S., Kurten, T., Zhang, Z. F., Gold, A., Surratt, J. D., Shilling, J. E., and Thornton, J. A.:
1001 Efficient Isoprene Secondary Organic Aerosol Formation from a Non-IEPDX Pathway, *Environ Sci*
1002 *Technol*, 50, 9872-9880, <https://doi.org/10.1021/acs.est.6b01872>, 2016.

1003 Liu, S., Shilling, J. E., Song, C., Hiranuma, N., Zaveri, R. A., and Russell, L. M.: Hydrolysis of
1004 Organonitrate Functional Groups in Aerosol Particles, *Aerosol Science and Technology*, 46, 1359-
1005 1369, <https://doi.org/10.1080/02786826.2012.716175>, 2012.

1006 Liu, S., Russell, L. M., Sueper, D. T., and Onasch, T. B.: Organic particle types by single-particle
1007 measurements using a time-of-flight aerosol mass spectrometer coupled with a light scattering
1008 module, *Atmospheric Measurement Techniques*, 6, 187-197, [https://doi.org/10.5194/amt-6-](https://doi.org/10.5194/amt-6-187-2013)
1009 [187-2013](https://doi.org/10.5194/amt-6-187-2013), 2013.

1010 Lopez-Hilfiker, F. D., Mohr, C., Ehn, M., Rubach, F., Kleist, E., Wildt, J., Mentel, T. F., Lutz, A.,
1011 Hallquist, M., Worsnop, D., and Thornton, J. A.: A novel method for online analysis of gas and
1012 particle composition: description and evaluation of a Filter Inlet for Gases and AEROSols
1013 (FIGAERO), *Atmospheric Measurement Techniques*, 7, 983-1001, [https://doi.org/10.5194/amt-](https://doi.org/10.5194/amt-7-983-2014)
1014 [7-983-2014](https://doi.org/10.5194/amt-7-983-2014), 2014.

1015 Lopez-Hilfiker, F. D., Mohr, C., Ehn, M., Rubach, F., Kleist, E., Wildt, J., Mentel, T. F., Carrasquillo,
1016 A. J., Daumit, K. E., Hunter, J. F., Kroll, J. H., Worsnop, D. R., and Thornton, J. A.: Phase partitioning
1017 and volatility of secondary organic aerosol components formed from α -pinene ozonolysis and OH
1018 oxidation: the importance of accretion products and other low volatility compounds,
1019 *Atmospheric Chemistry and Physics*, 15, 7765-7776, [https://doi.org/10.5194/acp-15-7765-](https://doi.org/10.5194/acp-15-7765-2015)
1020 [2015](https://doi.org/10.5194/acp-15-7765-2015).

1021 Lopez-Hilfiker, F. D., Mohr, C., D'Ambro, E. L., Lutz, A., Riedel, T. P., Gaston, C. J., Iyer, S., Zhang,
1022 Z., Gold, A., Surratt, J. D., Lee, B. H., Kurten, T., Hu, W. W., Jimenez, J., Hallquist, M., and Thornton,
1023 J. A.: Molecular Composition and Volatility of Organic Aerosol in the Southeastern U.S.:
1024 Implications for IEPOX Derived SOA, *Environ Sci Technol*, 50, 2200-2209,
1025 <https://doi.org/10.1021/acs.est.5b04769>, 2016.

1026 Mohr, C., Lopez-Hilfiker, F. D., Yli-Juuti, T., Heitto, A., Lutz, A., Hallquist, M., D'Ambro, E. L.,
1027 Rissanen, M. P., Hao, L. Q., Schobesberger, S., Kulmala, M., Mauldin, R. L., Makkonen, U., Sipila,
1028 M., Petaja, T., and Thornton, J. A.: Ambient observations of dimers from terpene oxidation in the
1029 gas phase: Implications for new particle formation and growth, *Geophysical Research Letters*, 44,
1030 2958-2966, <https://doi.org/10.1002/2017gl072718>, 2017.

1031 Murphy, D. M., Middlebrook, A. M., and Warshawsky, M.: Cluster analysis of data from the
1032 Particle Analysis by Laser Mass Spectrometry (PALMS) instrument, *Aerosol Science and*
1033 *Technology*, 37, 382-391, <https://doi.org/10.1080/02786820300971>, 2003.

1034 Pinero-Garcia, F., Ferro-Garcia, M. A., Chham, E., Cobos-Diaz, M., and Gonzalez-Rodelas, P.: A
1035 cluster analysis of back trajectories to study the behaviour of radioactive aerosols in the south-
1036 east of Spain, *J Environ Radioactiv*, 147, 142-152, <https://doi.org/10.1016/j.jenvrad.2015.05.029>,
1037 2015.

1038 Praske, E., Otkjaer, R. V., Crouse, J. D., Hethcox, J. C., Stoltz, B. M., Kjaergaard, H. G., and
1039 Wennberg, P. O.: Atmospheric autoxidation is increasingly important in urban and suburban
1040 North America, *P Natl Acad Sci USA*, 115, 64-69, <https://doi.org/10.1073/pnas.1715540115>, 2018.

1041 Rebotier, T. P., and Prather, K. A.: Aerosol time-of-flight mass spectrometry data analysis: A
1042 benchmark of clustering algorithms, *Anal Chim Acta*, 585, 38-54,
1043 <https://doi.org/10.1016/j.aca.2006.12.009>, 2007.

1044 Reitz, P., Zorn, S. R., Trimborn, S. H., and Trimborn, A. M.: A new, powerful technique to analyze
1045 single particle aerosol mass spectra using a combination of OPTICS and the fuzzy c-means
1046 algorithm, *Journal of Aerosol Science*, 98, 1-14, <https://doi.org/10.1016/j.jaerosci.2016.04.003>,
1047 2016.

1048 Roth, A., Schneider, J., Klimach, T., Mertes, S., van Pinxteren, D., Herrmann, H., and Borrmann, S.:
1049 Aerosol properties, source identification, and cloud processing in orographic clouds measured by
1050 single particle mass spectrometry on a central European mountain site during HCCT-2010,
1051 *Atmospheric Chemistry and Physics*, 16, 505-524, <https://doi.org/10.5194/acp-16-505-2016>,
1052 2016.

1053 Schobesberger, S., D'Ambro, E. L., Lopez-Hilfiker, F. D., Mohr, C., and Thornton, J. A.: A model
1054 framework to retrieve thermodynamic and kinetic properties of organic aerosol from
1055 composition-resolved thermal desorption measurements, *Atmospheric Chemistry and Physics*,
1056 18, 14757-14785, <https://doi.org/10.5194/acp-18-14757-2018>, 2018.

1057 Song, X. H., Hopke, P. K., Ferguson, D. P., and Prather, K. A.: Classification of single particles
1058 analyzed by ATOFMS using an artificial neural network, *ART-2A, Anal Chem*, 71, 860-865,
1059 <https://doi.org/10.1021/ac9809682>, 1999.

1060 Stolzenburg, D., Fischer, L., Vogel, A. L., Heinritzi, M., Schervish, M., Simon, M., Wagner, A. C.,
1061 Dada, L., Ahonen, L. R., Amorim, A., Baccarini, A., Bauer, P. S., Baumgartner, B., Bergen, A., Bianchi,
1062 F., Breitenlechner, M., Brilke, S., Mazon, S. B., Chen, D. X., Dias, A., Draper, D. C., Duplissy, J.,
1063 Haddad, I., Finkenzeller, H., Frege, C., Fuchs, C., Garmash, O., Gordon, H., He, X., Helm, J.,
1064 Hofbauer, V., Hoyle, C. R., Kim, C., Kirkby, J., Kontkanen, J., Kuerten, A., Lampilahti, J., Lawler, M.,
1065 Lehtipalo, K., Leiminger, M., Mai, H., Mathot, S., Mentler, B., Molteni, U., Nie, W., Nieminen, T.,
1066 Nowak, J. B., Ojdanic, A., Onnela, A., Passananti, M., Petaja, T., Quelever, L. L. J., Rissanen, M. P.,
1067 Sarnela, N., Schallhart, S., Tauber, C., Tome, A., Wagner, R., Wang, M., Weitz, L., Wimmer, D.,
1068 Xiao, M., Yan, C., Ye, P., Zha, Q., Baltensperger, U., Curtius, J., Dommen, J., Flagan, R. C., Kulmala,
1069 M., Smith, J. N., Worsnop, D. R., Hansel, A., Donahue, N. M., and Winkler, P. M.: Rapid growth of
1070 organic aerosol nanoparticles over a wide tropospheric temperature range, *P Natl Acad Sci USA*,
1071 115, 9122-9127, <https://doi.org/10.1073/pnas.1807604115>, 2018.

1072 Takahama, S., Gilardoni, S., Russell, L. M., and Kilcoyne, A. L. D.: Classification of multiple types
1073 of organic carbon composition in atmospheric particles by scanning transmission X-ray
1074 microscopy analysis, *Atmos Environ*, 41, 9435-9451,
1075 <https://doi.org/10.1016/j.atmosenv.2007.08.051>, 2007.

1076 Wang, D. S., and Ruiz, L. H.: Chlorine-initiated oxidation of n-alkanes under high-NO_x conditions:
1077 insights into secondary organic aerosol composition and volatility using a FIGAERO-CIMS,
1078 *Atmospheric Chemistry and Physics*, 18, 15535-15553, [https://doi.org/10.5194/acp-18-15535-](https://doi.org/10.5194/acp-18-15535-2018)
1079 [2018](https://doi.org/10.5194/acp-18-15535-2018), 2018.

1080 Wegner, T., Hussein, T., Hameri, K., Vesala, T., Kulmala, M., and Weber, S.: Properties of aerosol
1081 signature size distributions in the urban environment as derived by cluster analysis, *Atmos*
1082 *Environ*, 61, 350-360, <https://doi.org/10.1016/j.atmosenv.2012.07.048>, 2012.

1083 Zhao, W. X., Hopke, P. K., and Prather, K. A.: Comparison of two cluster analysis methods using
1084 single particle mass spectra, Atmos Environ, 42, 881-892,
1085 <https://doi.org/10.1016/j.atmosenv.2007.10.024>, 2008.

1086 Zhao, Y., Thornton, J. A., and Pye, H. O. T.: Quantitative constraints on autoxidation and dimer
1087 formation from direct probing of monoterpene-derived peroxy radical chemistry, P Natl Acad Sci
1088 USA, 115, 12142-12147, <https://doi.org/10.1073/pnas.1812147115>, 2018.

1089 Zhou, L. M., Hopke, P. K., and Venkatachari, P.: Cluster analysis of single particle mass spectra
1090 measured at Flushing, NY, Anal Chim Acta, 555, 47-56, <https://doi.org/10.1016/j.aca.2005.08.061>,
1091 2006.

1092

10. Tables

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

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

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

Expt #	SOA type	Pre-processing						Clustering				
		N_{total}	$N_{\text{anomalous}}$	N_{filtered}	$f_{\text{m,filtered}}$	ζ_{ref}	$f_{\text{m,ref}}$	ε	N_{c}	$N_{\text{c,one}}$	$f_{\text{m,unclustered}}$	$R_{\text{interClst}}$
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

$N_{\text{anomalous}}$ – Number of anomalous ions

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

$f_{\text{m,filtered}}$ – 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

$f_{\text{m,ref}}$ – 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

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

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

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

11. Figures

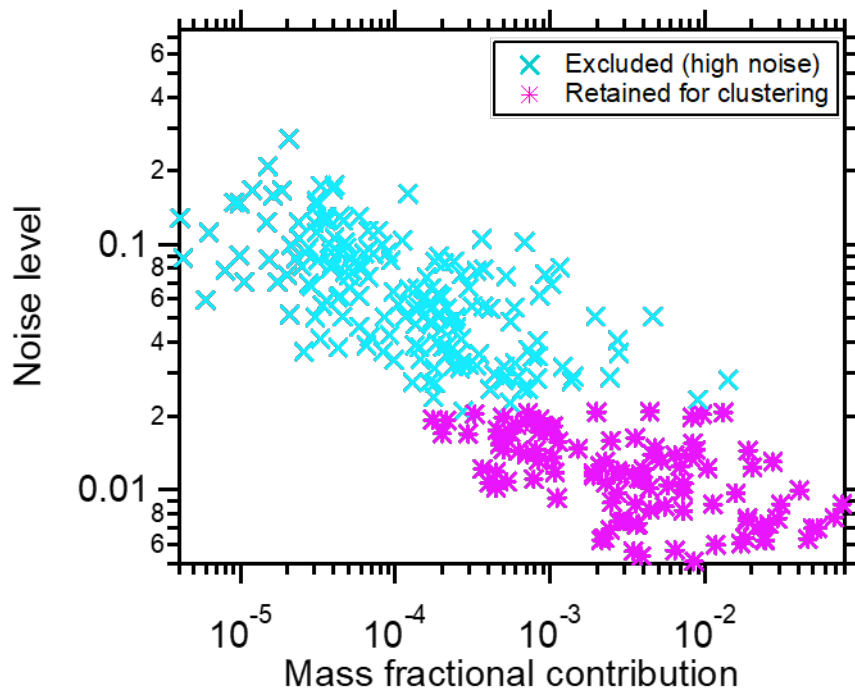


Figure 1: The relationship between thermogram noise levels and the fractional contributions of the corresponding ions to total mass, for α -pinene + OH SOA. The noise threshold, $\xi_{\text{ref}} = 0.021$ and is used to distinguish high-noise thermograms (cyan markers) from thermograms having acceptable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

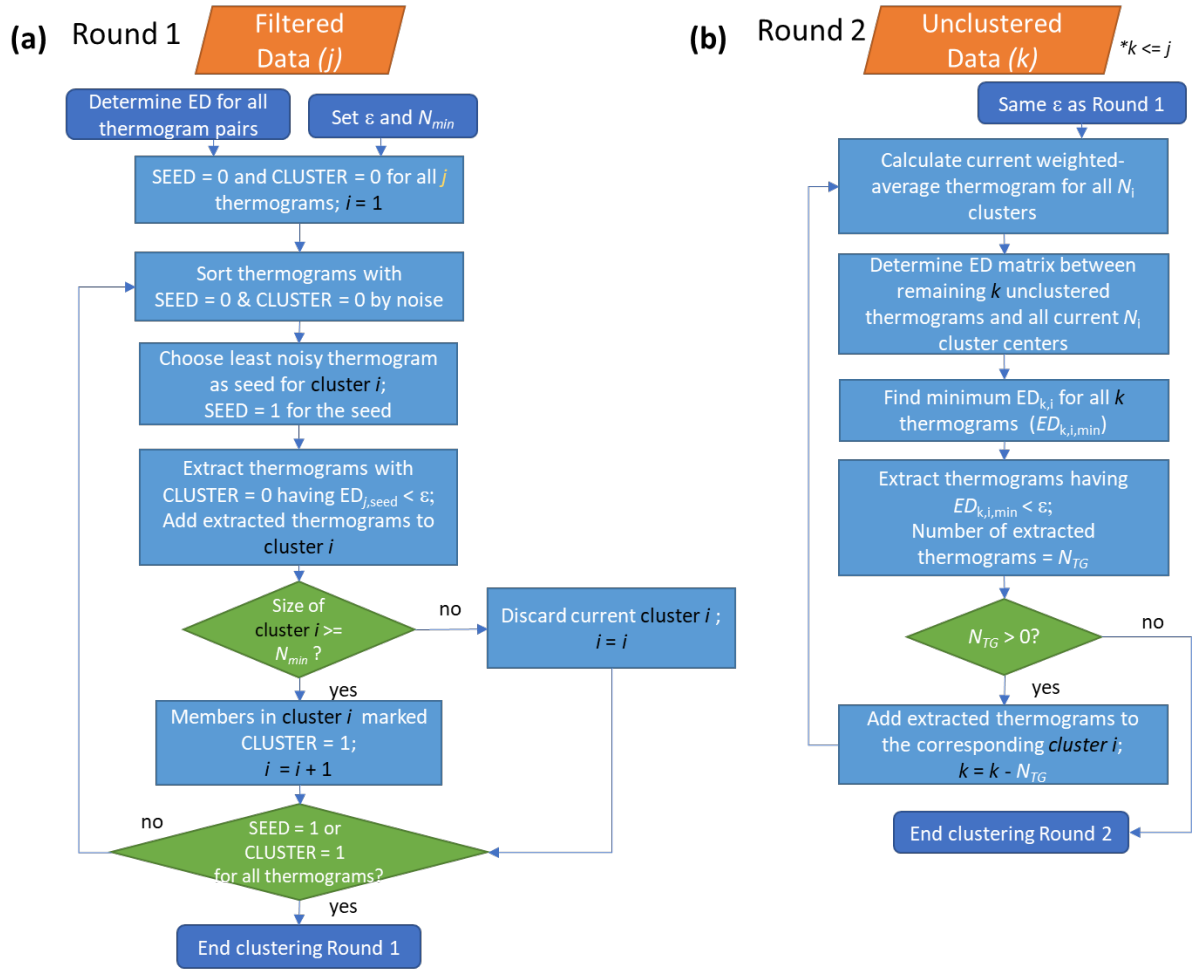


Figure 2: Flow of the noise-sorted scanning clustering. There are two rounds of clustering. (a) Round 1: The ED between all thermogram pairs are calculated and two parameters, ϵ and N_{min} , are set. Each thermogram is initialized with state SEED = 0 and CLUSTER = 0. Only thermograms with SEED = 0 and 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.

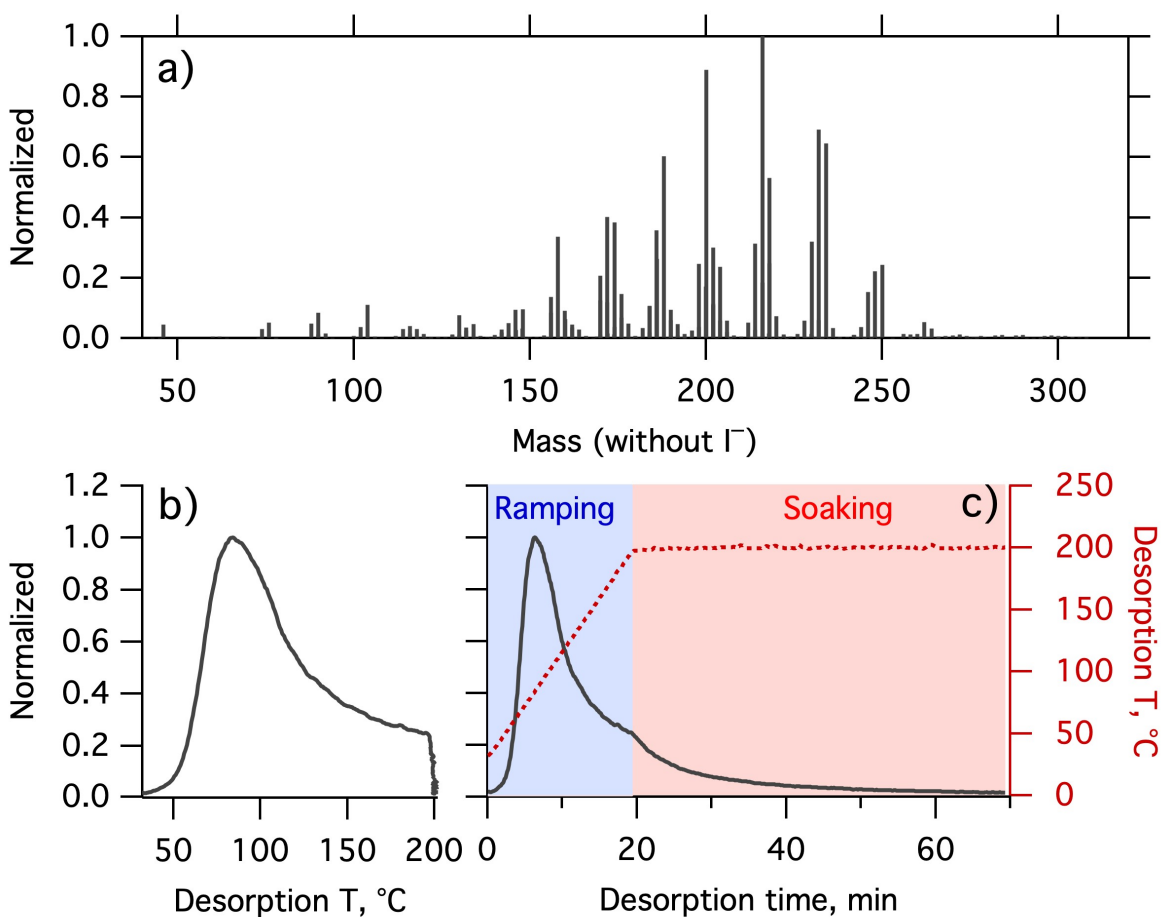


Figure 3. (a) Mass spectrum of α -pinene + OH SOA measured by FIGAERO-CIMS. The mass excludes iodine. (b) Normalized thermogram of the bulk SOA versus temperature. (c) Normalized thermogram of the bulk SOA versus time (black line) and the variation in desorption temperature with time (dark red dashed line). The long tail during the soaking period is evident when the thermogram is considered in time space. The light blue shaded area denotes the ramping period and the pink shaded area the soaking period.

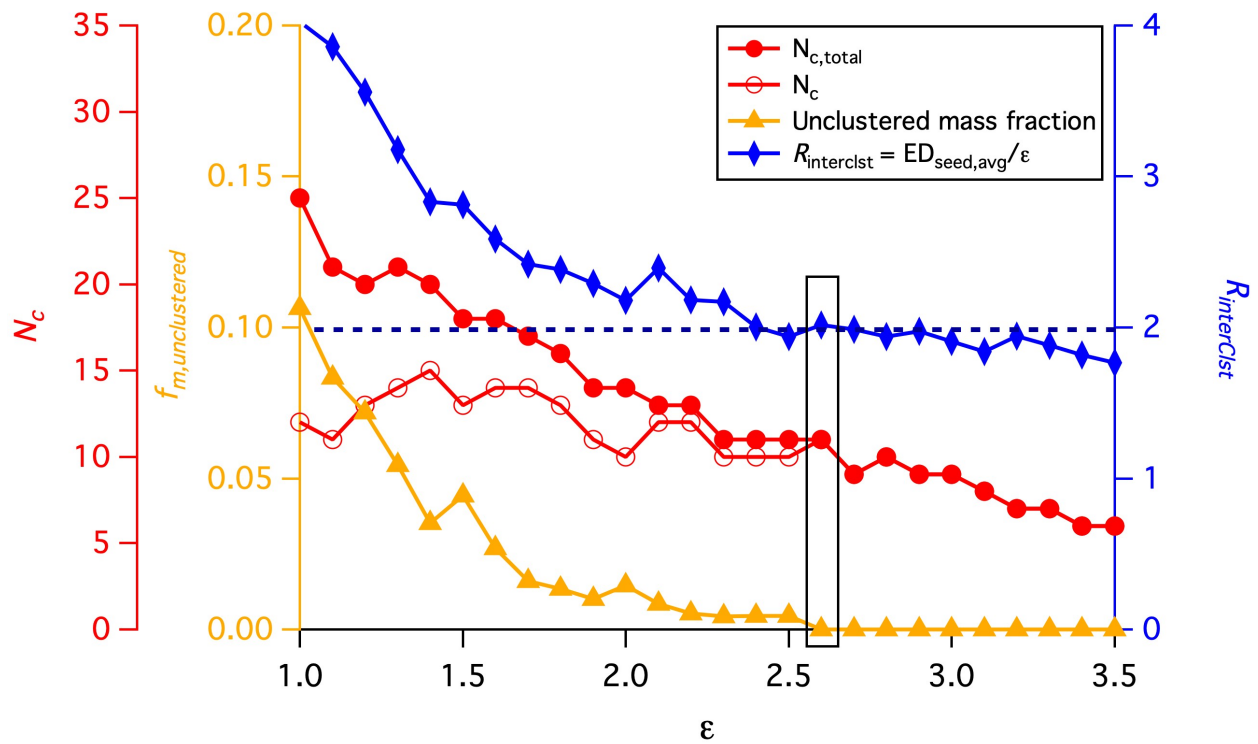


Figure 4. The variation of four parameters, N_c , $N_{c,total}$, $f_{m,unclustered}$ and $R_{interClst}$ as a function of the distance criterion ε . The black horizontal dashed line guides the judgement for $R_{interClst} \geq 2$. The values highlighted by a rectangle are the values corresponding to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

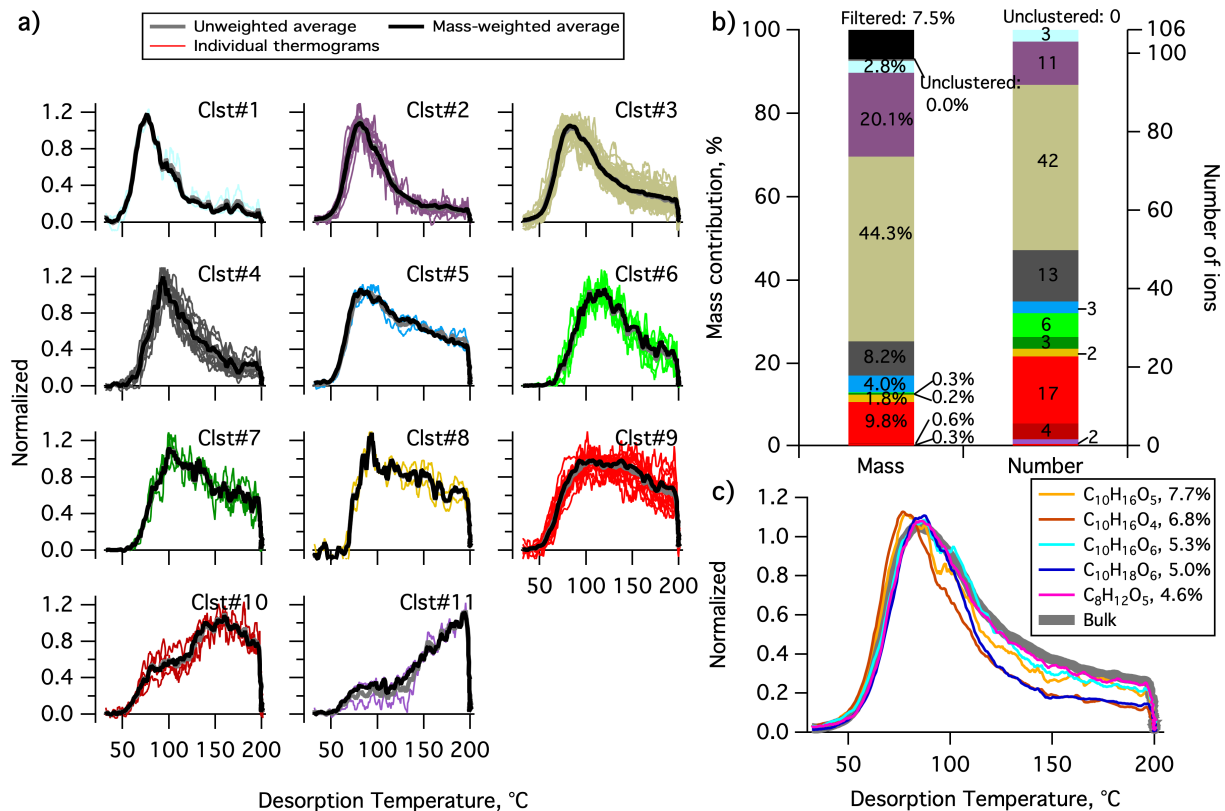


Figure 5. Clustering results for α -pinene + OH SOA. (a) Unweighted average thermograms (bold grey lines), mass-weighted average thermograms (bold black lines) and individual members (colored lines) of the 11 clusters identified. (b) Percentage contribution of each cluster to the total mass, as well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

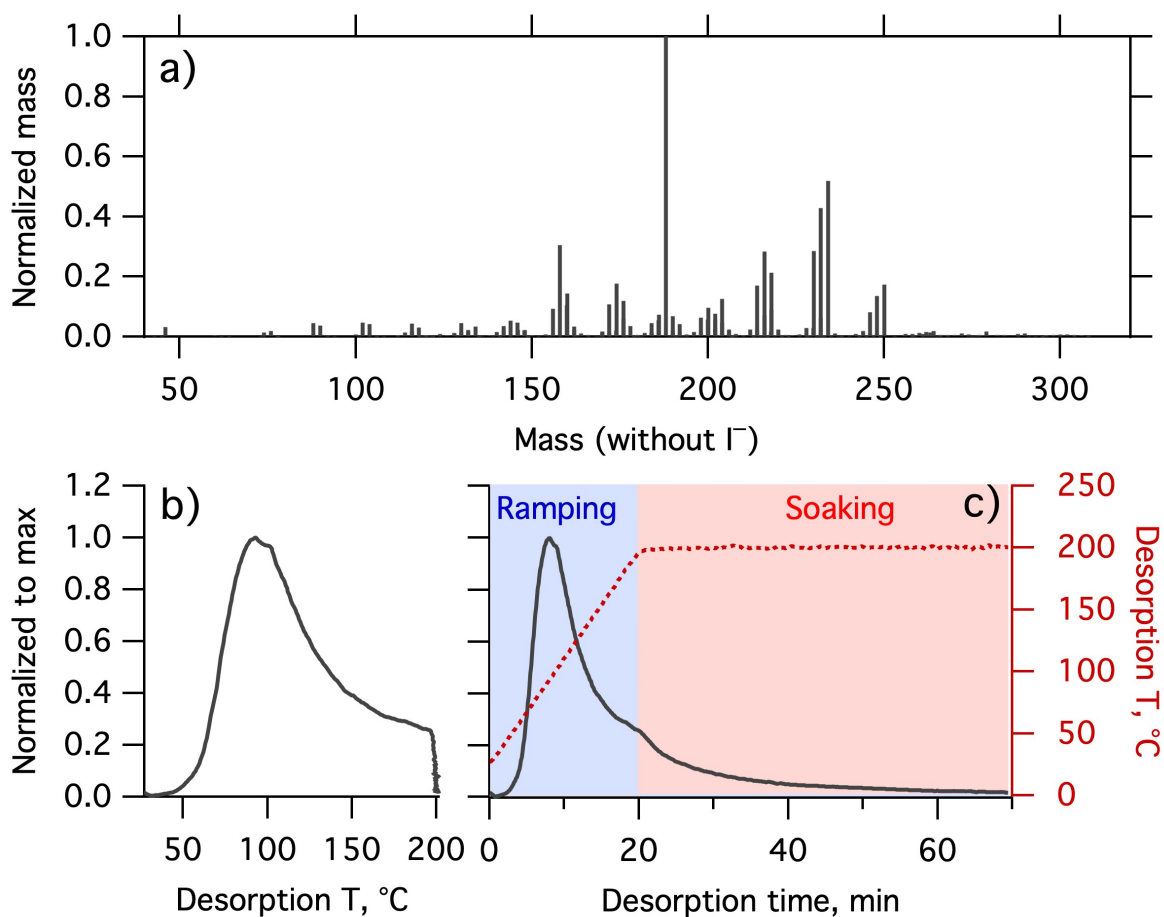


Figure 6. Same as Figure 3, but for Δ -3-carene + OH SOA. (a) SOA mass spectrum measured by FIGAERO-CIMS. The mass excludes iodine. The normalized thermogram of the bulk SOA versus (b) temperature and (c) time. In (c) the light blue shaded area denotes the ramping period and the pink shaded area the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

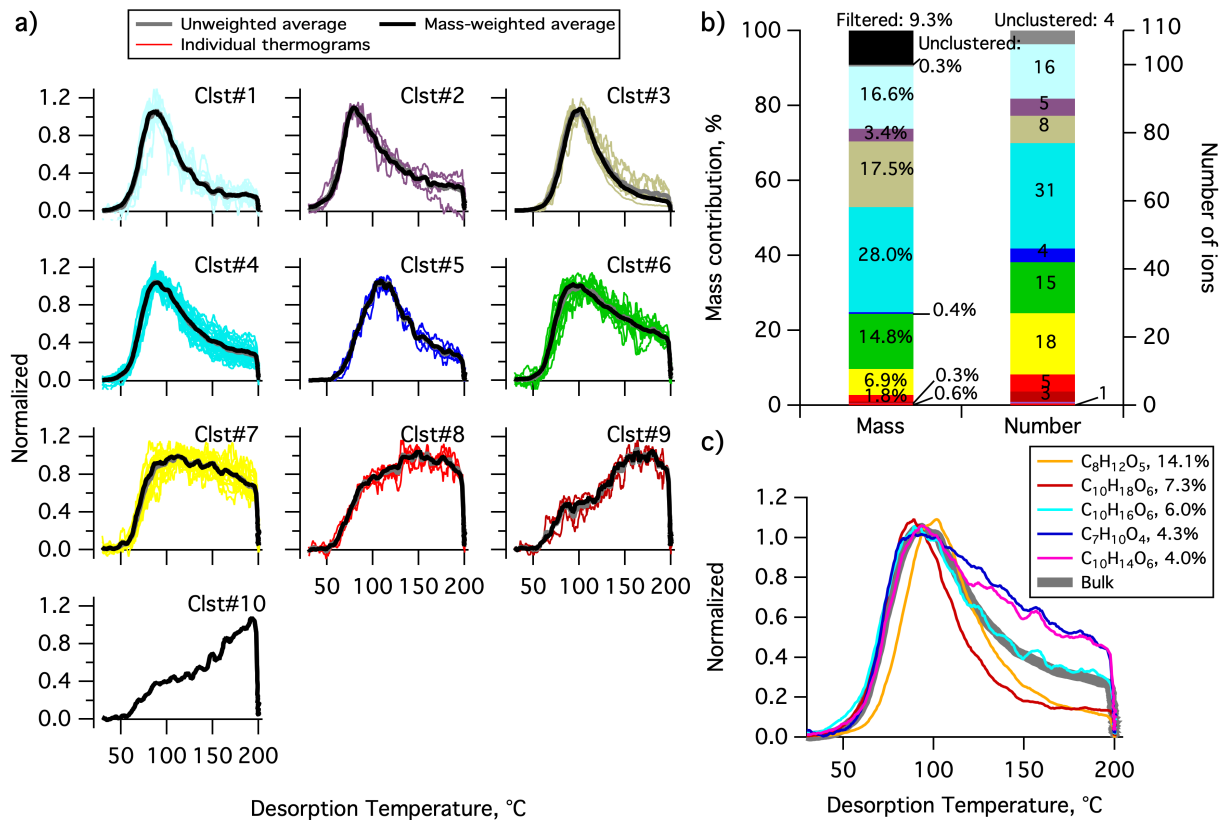


Figure 7. Same as Figure 5, but for Δ -3-carene + OH SOA. (a) Unweighted average thermograms (bold grey lines), mass-weighted average thermograms (bold black lines) and individual members (colored lines) of the ten clusters identified. (b) Percentage contribution of each cluster to the total mass, as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

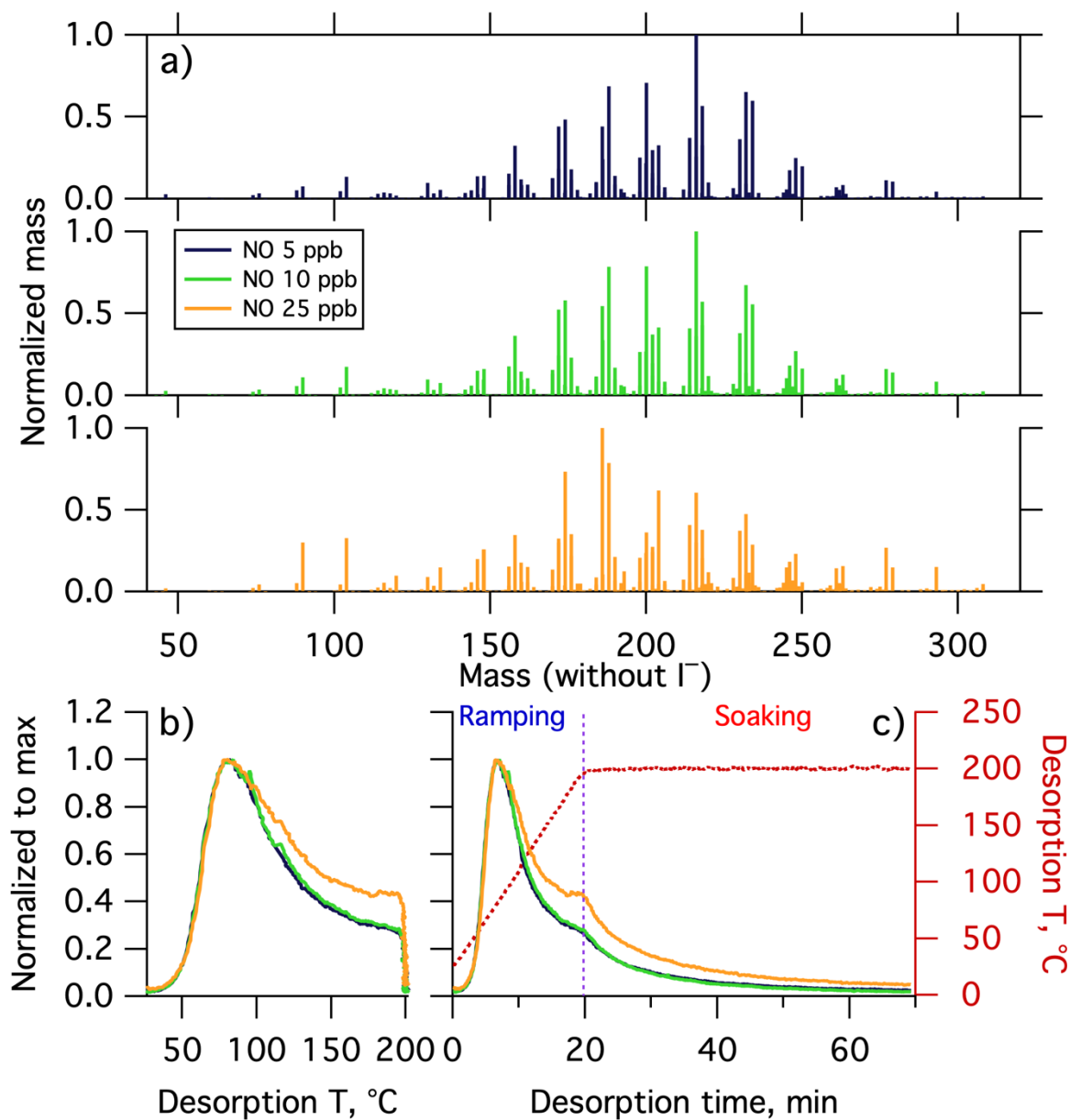


Figure 8. (a) Mass spectra of α -pinene + OH SOA formed with different NO concentrations, normalized to the most abundant ions mass concentration. The mass excludes iodine. Normalized thermograms of the bulk SOA versus (b) temperature and (c) desorption time, with the desorption temperature shown in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

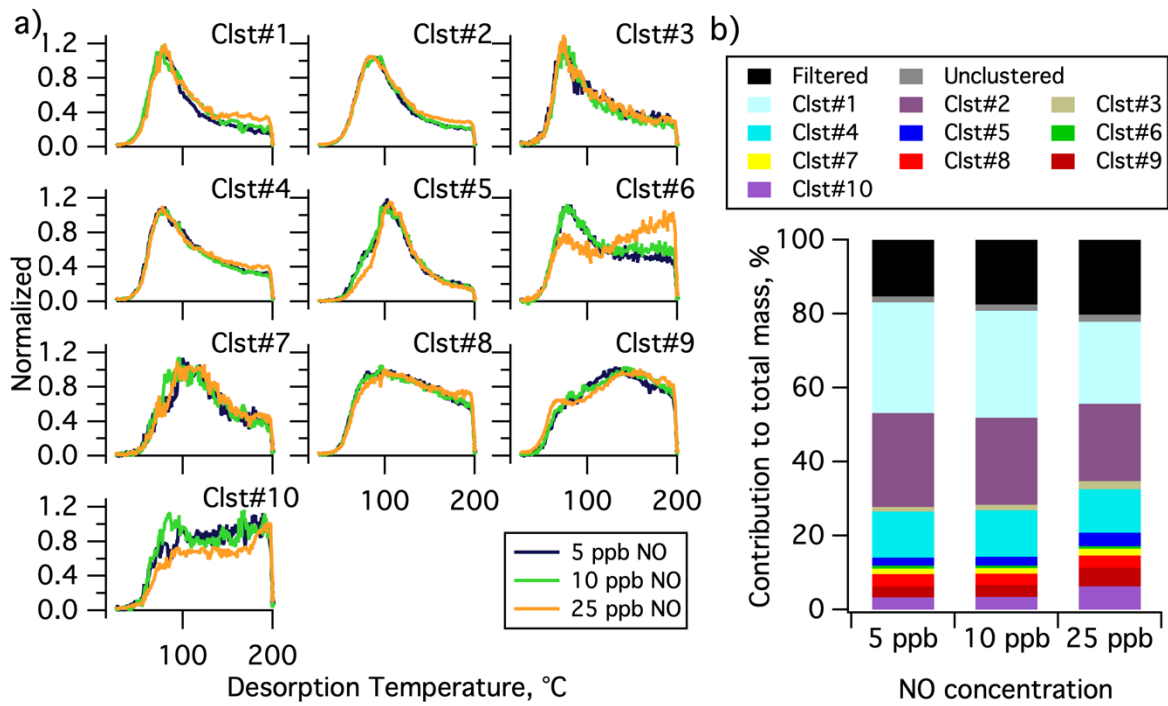


Figure 9. Single clustering results for α -pinene + OH SOA as a function of NO concentration. (a) Comparison of the normalized, weighted average thermograms of the ten clusters for the 5 ppb NO (navy), 10 ppb NO (green) and 25 ppb NO (orange) experiments. (b) Contribution of each cluster to the total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

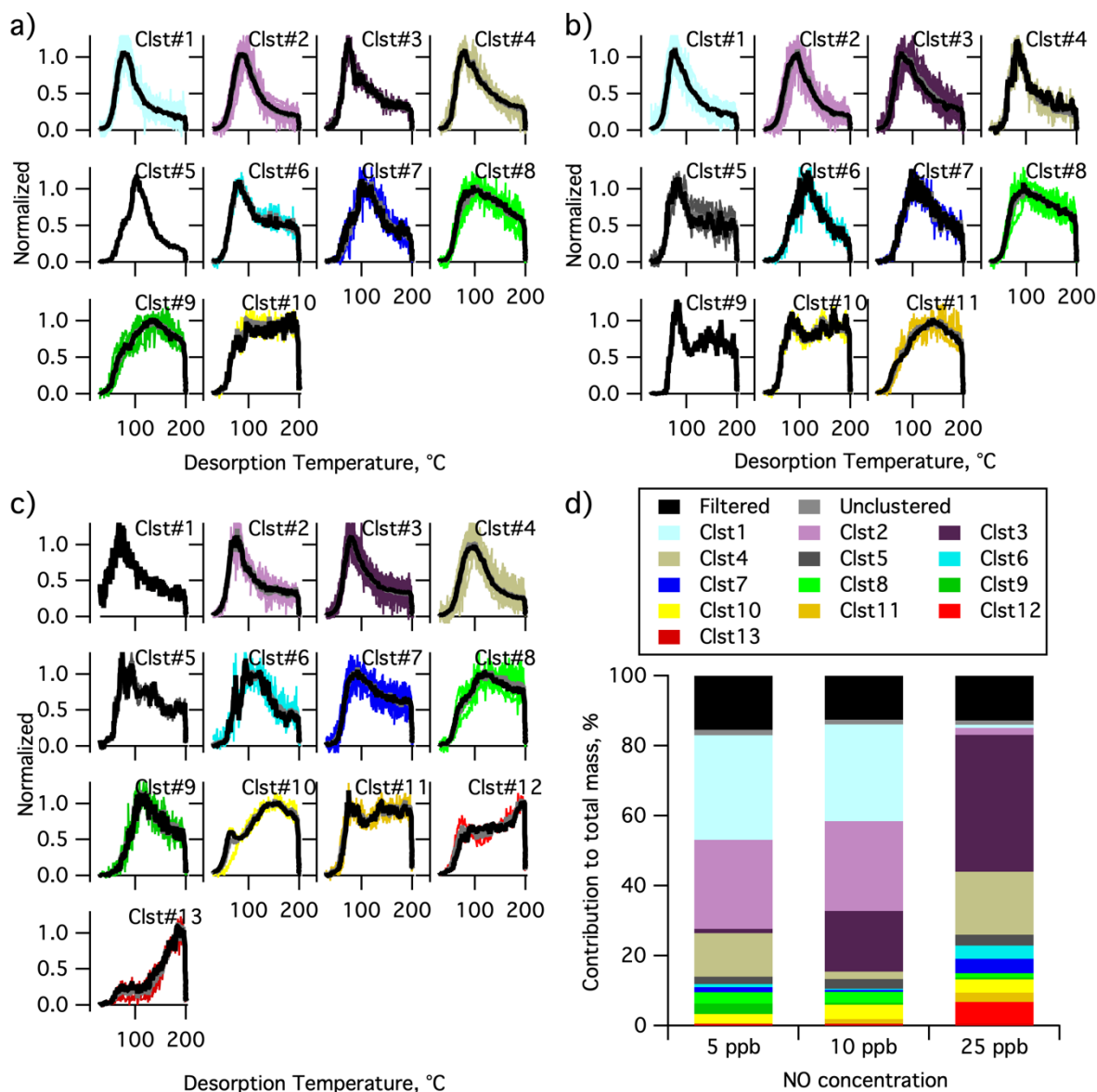


Figure 10. Multiple clustering results for α -pinene + OH SOA as a function of NO concentration. Clustering results are separately shown for the (a) 5 ppb NO, (b) 10 ppb NO, and (c) 25 ppb NO experiments. Each panel includes unweighted average thermograms (grey lines), mass-weighted average thermograms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

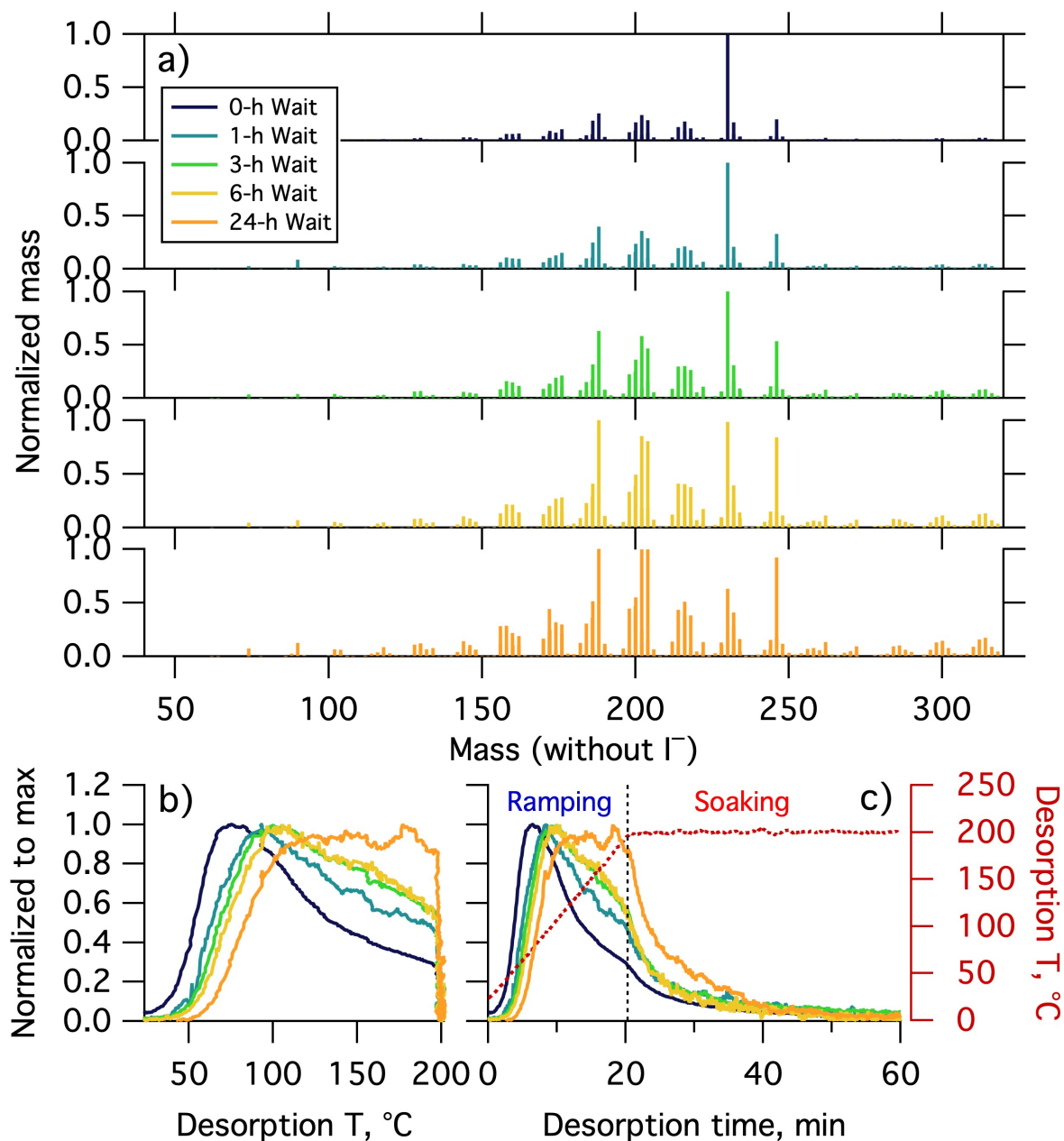


Figure 11. (a) Normalized mass spectra of α -pinene + O_3 SOA measured after different extents of isothermal evaporation at room temperature. The mass excludes iodine. The normalized thermograms of bulk SOA versus (b) temperature and (c) time, with the desorption temperature shown as a red dashed line. The vertical black dashed line in (c) delineates between ramping and soaking. The mass spectrum or thermogram colors indicate the isothermal evaporation time (see legend), with darker colors indicating shorter times.

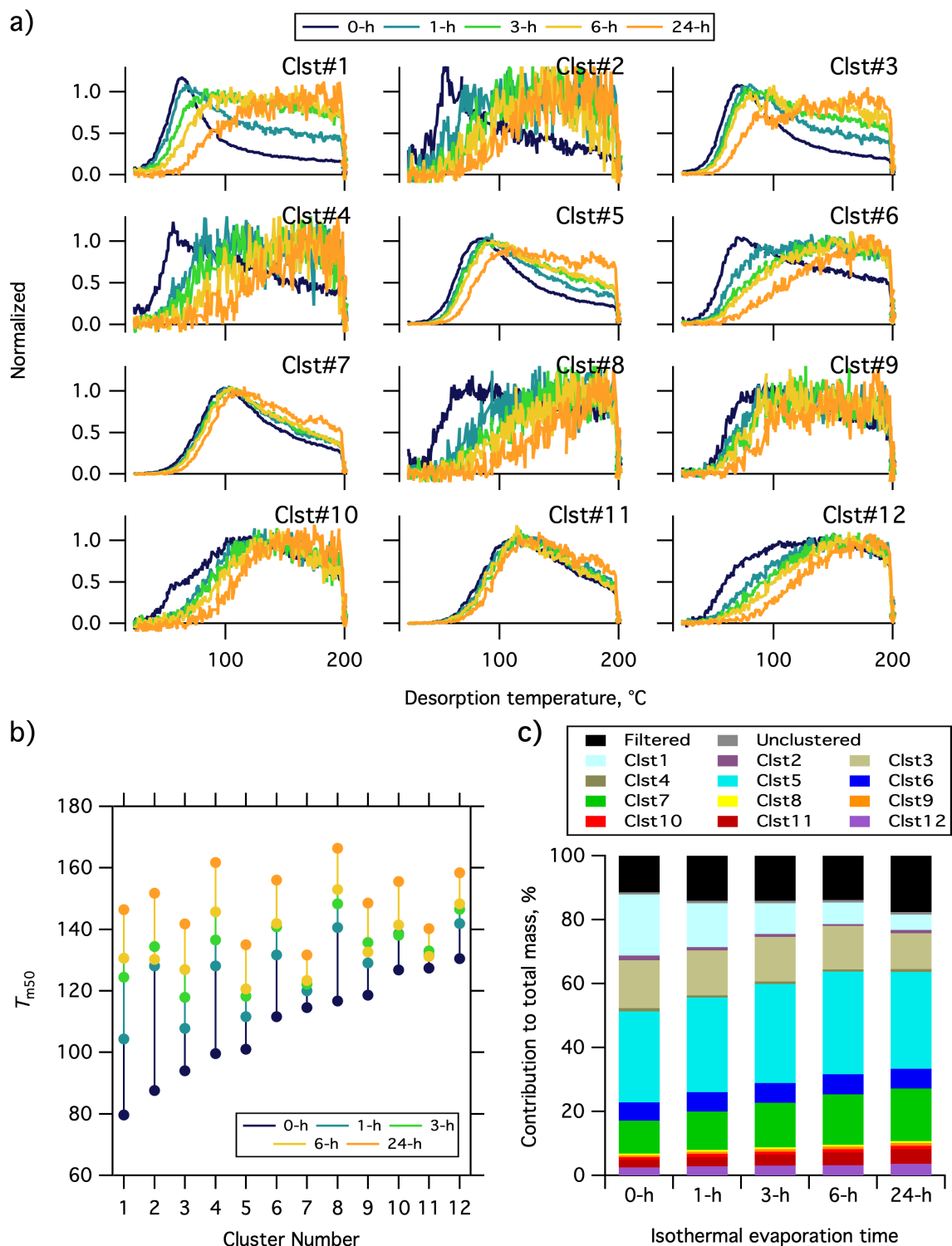


Figure 12. Single clustering results for α -pinene + O_3 SOA for different isothermal evaporation times. (a) Comparison of the normalized, weighted-average thermograms of the 12 clusters of 0-h wait (navy), 1-h wait (blue), 3-h wait (green), 6-h wait (yellow) and 24-h wait (orange) experiments. Note that the absolute signals of all of the clusters decrease with evaporation, but to varying extents (**Error! Reference source not found.**).

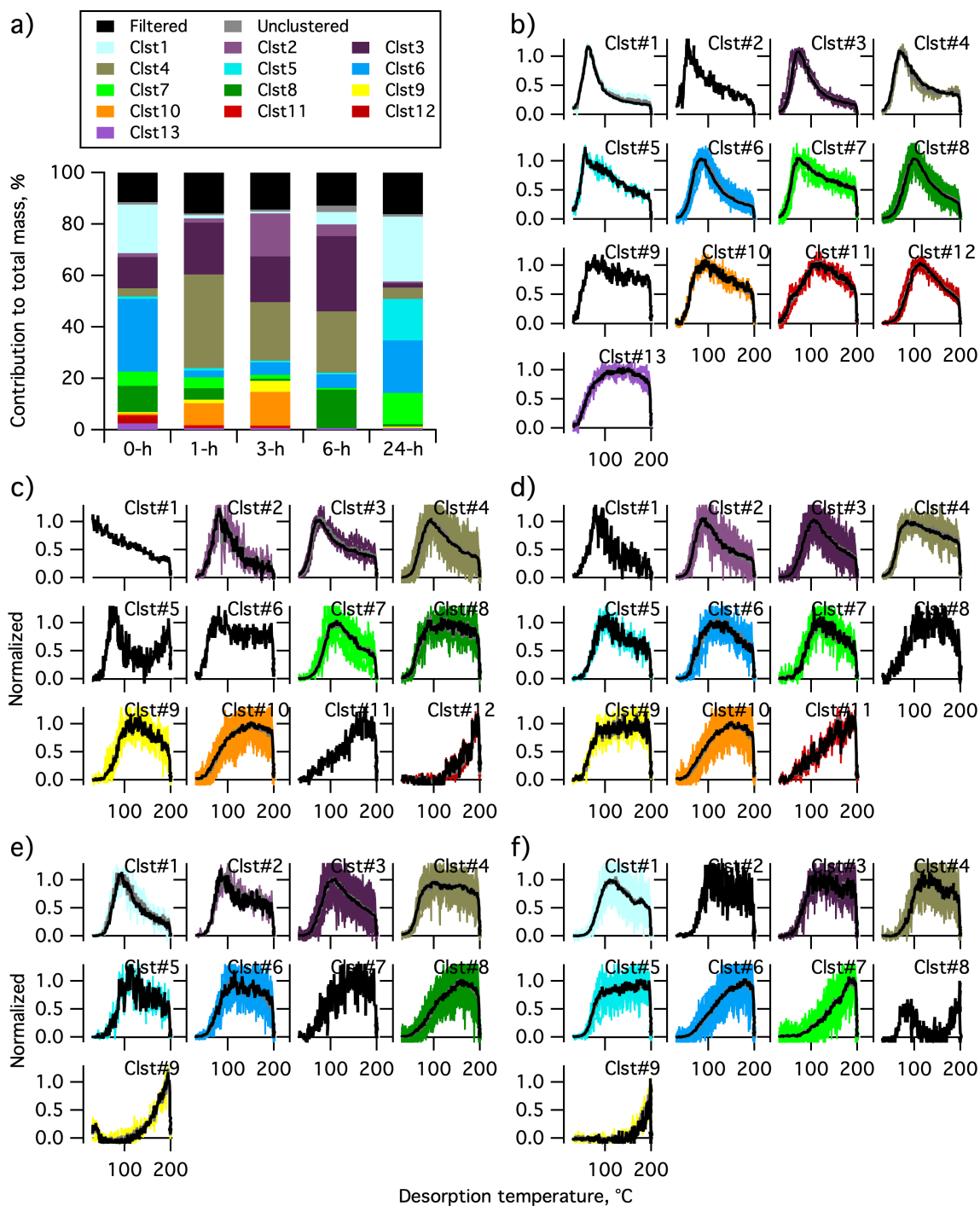


Figure 13. Multiple clustering results for α -pinene + O₃ SOA as a function of isothermal evaporation time. (a) Contribution of each cluster to the total mass for each experiment, along with the contributions of filtered-out ions (black bar) and unclustered ions (gray bar). The number of clusters obtained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

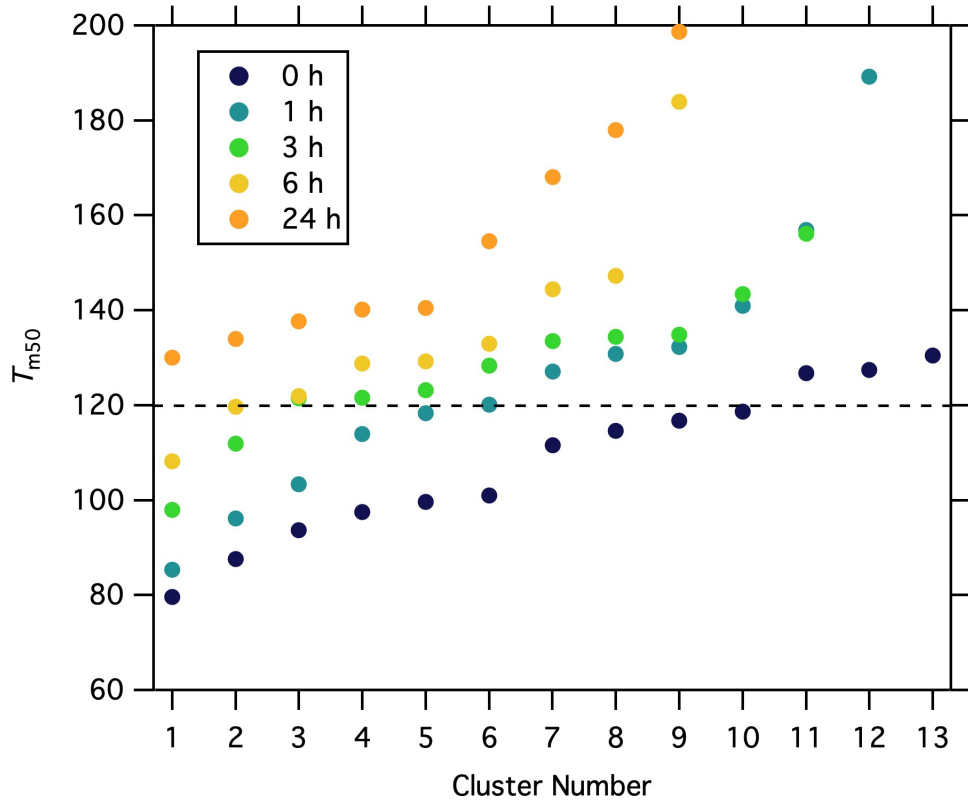


Figure 14. The T_{m50} values of the cluster-specific thermograms from multiple clustering for the five isothermal evaporation experiments.