

We thank the reviewers for the thoughtful comments. We address each comment individually below, with the reviewers' initial comment in **black** and our responses in **blue**. A track changes version showing all changes made to the manuscript is appended at the end.

Response to Reviewer #1

The authors present a well structured detailed report of method they are proposing to analyse thermograms collected using FIGAERO-CIMS data. Although the manuscript focusses on cluster analysis, it clear that a considerable amount of work has gone into collecting the data, and developing and trialling the method (Noise-Sorted Scanning Clustering). It is difficult to find fault in the work. Their introduction gives a good panorama of the cluster analysis and air quality data vista. They select the various suitable cluster analysis methods and make comparisons using their data before justifying their choice of NSSC. The data flow is well described and supported by illustrations and them exemplified by application to laboratory generated SOA. It will be interesting to see how this method deals with ambient data. The following points are simply minor considerations on how to improve the presentation taking into consideration this is a paper on a new method of clustering and as a reader, I am asking if I could reproduce this method for a different application or in a different code.

We first thank the reviewer for the positive assessment. We think the NSSC method is fairly easy to transfer to a different code. The application of NSSC can be potentially expanded to any composition-resolved data sets, such as diurnal changes of different compounds measured in ambient air, temporal changes of different generations of species in a smog chamber, and composition-dependent size distributions. All of the above data sets share a common property that the noise of the curve/spectrum is related to the composition. We have expanded the discussion of the application of NSSC in section 5.

“This paper focuses only on the description of the clustering algorithm and its potential as a tool to characterize the thermal properties of organic aerosol in further detail. The application of NSSC can be potentially expanded to any other composition-resolved data sets, such as diurnal changes of different compounds measured in ambient air, temporal changes of different generations of species in a smog chamber, and composition-dependent size distributions. All of the above data sets share a common property that the noise of the curve/spectrum is related to the composition. Therefore, NSSC would facilitate the analysis by taking noise into consideration.”

To the reviewers point that this will be interesting to see how the method deals with ambient data, we agree and are actively pursuing this idea.

1. I appreciate the descriptions given of the methods and especially figure 2 and I am asking myself if more detailed mathematics be included to describe the method?

The reviewer raises an important point. Figure 2 serves to provide an overview of the flow of NSSC method to help readers understand the algorithm. Therefore, we made figure 2 a more generic description. We described all of the detailed mathematics in the text only because most of these parameters are ultimately data-specific and user-defined.

2. I can see that there is a lot of information conveyed in figures 5, 7, 9, 10 and especially 13 and I am asking myself if they can be enhanced to better convey their message?

The reviewer raises an important point. There are indeed a lot of information in each of the figure the reviewer mentions. For figure 5 and 7, we think they are the best way to provide clustering results for simple single-precursor SOA systems at the moment. We have tried to plot all the average cluster thermograms in one figure. However, due to the existence of many overlaps, it makes the comparison between different clusters more difficult. We have also tried to use pie chart instead of bars to show the percentage contribution of clusters. We think they are equally efficient in conveying information and the bar chart is more convenient when different systems are compared in one figure. For figure 10 and 13, they show the clustering results of a set of experiments using the multiple clustering approach. It is necessary to show all the averaged thermograms of all the experiments in order to find which of the thermograms is common in different experiments while which of the thermogram is unique in only one experiment. As for the bar charts, it would be clearer to show the percentage contribution of grouped clusters based on T_{m50} as is described in section 4.4.2. However, we choose to leave the original, detailed information in the figure to both give an example of the detailed clustering results of NSSC and let the readers explore further interpretations of the clustering results. Therefore, we have not made modifications to these figures. We note that we have provided the information from these figures in an accessible, downloadable format with the associated dataset so that readers can explore the data further.

Response to Reviewer #2

Note: This review is focused on the clustering aspects of the paper only. I make no comment on whether the generated clusters constitute a useful and valid contribution to the analysis of composition-dependent organic aerosol thermal desorption measurements as this is not in my area of expertise. I leave that for others to comment on.

From the title and abstract I was expecting to read about a new clustering algorithm. However, this paper presents a data pre-processing step before using DBScan, followed by an application-specific post-clustering process.

Summary The section on the clustering process should be re-written. This is not a new clustering algorithm, but a method of pre- and post-processing data to provide useful analysis. I do not deny the usefulness of the analysis, its methodology or results. There are also points below that should be clarified.

We first thank the reviewer for the thoughtful comments on the clustering algorithm. The reviewer's comments have helped us to clarify aspects of the manuscript, making clearer the unique aspects of this work.

We agree that the NSSC stems from DBScan, and noted as much in the text when we state that the NSSC is "a novel variant of the DBSCAN algorithm." However, we contend that the NSSC differs from DBScan in important ways besides the process of seeds sorting and second-round clustering the reviewer mentions, making it sufficiently "new." For example, the way NSSC and DBScan define a cluster differ. Details will be described in the following responses. (Perhaps semantics, but we would also contend that the "preprocessing" and "post-clustering process" constitute part of the overall algorithm. A definition of algorithm is "a process or set of rules to be followed in calculations or other problem-solving operations, especially by a computer.") At the current stage, NSSC is designed specifically for FIGAERO-CIMS

thermograms, but has the potential to be applied to other composition-specific data sets. The values of many parameters and factors used in the data processing are empirically derived. Ultimately, they can be adjusted by users based on their specific applications. All this said, we have made revisions to the manuscript to more clearly indicate the link to the DBScan algorithm, including in the abstract.

1. DBScan is compared to other clustering algorithms, however these alternatives are significantly different in their intended use, i.e. they generate hyper-ellipsoidal clusters, whereas DBScan is specifically targeted at arbitrarily shaped clusters. It would be better to compare DBScan with similar algorithms. For the purposes of this paper, I am not sure that comparison to other algorithms is required.

The reviewer raises an important point about the comparison of different clustering algorithms. As the reviewer noted, DBScan and other algorithms have different intended use. For the FIGAERO-CIMS thermograms, however, it is difficult to define the data as hyper-ellipsoidal clusters or arbitrarily shaped clusters. Therefore, we tried DBScan, k-means, k-medoids and mean-shift, as they are well-known and most commonly used clustering algorithms. A brief description and comparison of these methods are presented in section 2.3 to give a context of why we chose NSSC. We believe this comparison provides value, especially as the only other attempt at clustering FIGAERO data that we are aware of used k-means.

2. Line 111 states 'a novel variant of DBScan'. The clustering algorithm used appears unchanged from DBScan, but rather a pre-processing technique of ordering the data is utilised before clustering.

The reviewer points out that the seed-sorting process of NSSC should not be considered as a variant of DBScan. We understand that this distinction is important to make, although contend that the "algorithm" encompasses the entire set of standardized procedures used, including pre-processing. By adding the pre-processing and post-processing steps, the NSSC is definitionally a "variant." (Variant, *noun*, "a form or version of something that differs in some respect from other forms of the same thing or from a standard.") We also note that there is an additional aspect that makes NSSC different from DBScan. To the best of our knowledge, in a cluster defined by DBScan, there are a core point, directly reachable points and reachable points. Directly reachable points are within the critical distance ϵ of the core point, while reachable points are within the distance ϵ of directly reachable points or other reachable points. The inclusion of reachable points is the key reason why DBScan can find arbitrarily shaped clusters. However, NSSC only considers the core point (seed) and the directly reachable points (neighbors) as a cluster in the first step. There is a second step of clustering where the seed is redefined based on the cluster average thus far and new directly reachable points added, expanding the number of members that are included in a given cluster. However, the number of new members added in this second step tends to be small. In some ways, NSSC is more similar to for example k-means in a way it generates hyper-ellipsoidal clusters. We have clarified the difference between NSSC and DBScan in section 2.3.

"Noisy members tend to naturally be excluded from any clusters. NSSC is a variant of DBSCAN. It does, however, differ from the standard DBSCAN algorithm because NSSC only searches for neighbors of the seed, while DBSCAN also searches for neighbors of the neighbors. As such, the sorting of seeds by noise levels is a key aspect of the NSSC algorithm which we have found provides for more robust clustering results."

3. Line 132 'absolute magnitude' is moot, magnitude has no direction.

We have deleted "absolute".

4. The paragraph from lines 134-144 describes a process based on a number of factors which are not justified. This leaves the questions: why 100 points; why +/- 50 points; and why is an anomaly only outside of -3σ and not outside of $+3\sigma$ also?

The reviewer raises an important point about justification of several values used in the pre-processing analysis. These values are derived based on consideration of 10 different smog chamber experiments with different chemical systems but similar FIGAERO-CIMS operation. The number of points we determined should be used for noise determination (100 points) derives from inspection of the thermograms from the different experiments and our finding that, as stated, during this period the “signals are usually relatively constant.” Use of more points leads to undesirable incorporation of times when the signals are still declining during the soaking period, increasing the standard deviation. Use of fewer points leads to larger overall noise levels. We have added statements to this effect to the manuscript. The number of points used to determine the minimum signal (+/-50) was determined based on the temperature ramping speed and a desire to identify as “noisy” only those thermograms that exhibited large negative deviations. We established that use of many fewer points led to an over-sensitivity to small fluctuations. Use of a greater number of points led to excessive smoothing. Additionally, the selection of +/- 50 points for calculating the minimum provides consistency with the number of points averaged for determining the noise. An anomaly refers to only outside of -3σ because the values in a thermogram are all background corrected and expected to be positive. So, $A_{\min} < -3\sigma$ indicates that the minimum is at least three standard deviations below zero. We have added a discussion of the values of factors used in this paper at the beginning of section 2.1.1 to clarify.

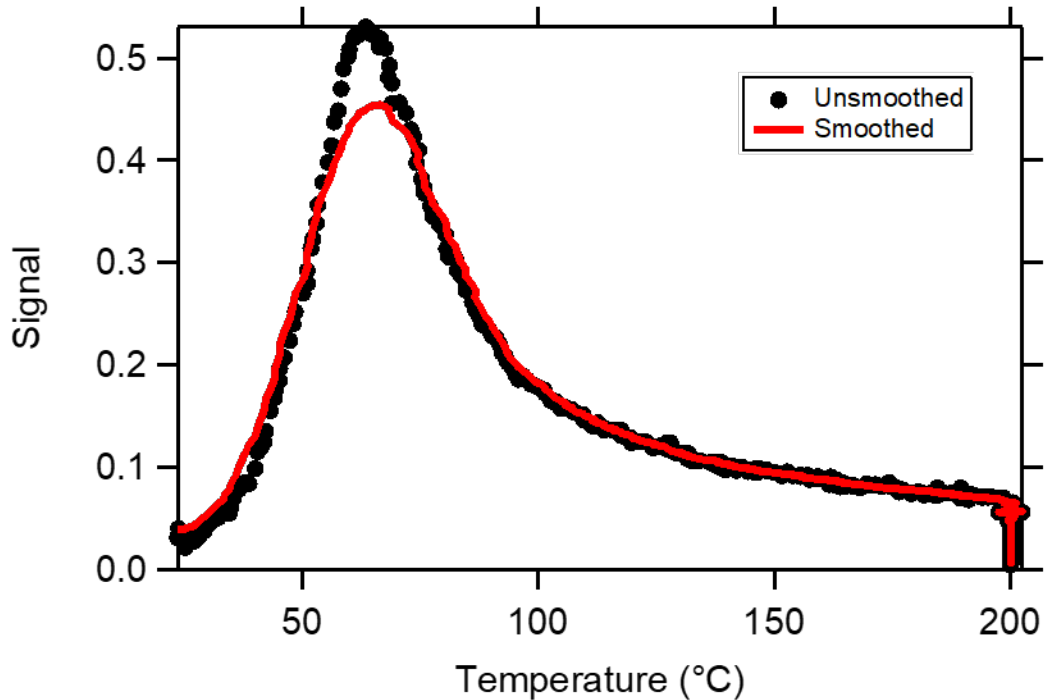
“Estimate a reference noise level (σ_{ref}) for each thermogram as the standard deviation of the last 100 points (corresponding to 500 seconds) of the thermogram at the end of the constant-temperature soaking period, during which the signals are usually relatively constant. *Use of more points incorporates times when the signals were still decreasing, while use of fewer points provides a less robust estimate of the noise level.* (ii) Find the minimum in the thermogram and calculate the average of this and the 50 points (corresponding to 250 seconds, or 100 points) before and after the minimum, A_{\min} . *This provides for consistency with the determination of σ_{ref}* (iii) Identify thermograms for which $A_{\min} < -3|\sigma_{\text{ref}}|$ as anomalous and exclude these associated ions from further analysis. In other words, when a thermogram has a valley with averaged negative values exceeding the magnitude of three times of the reference noise level, then it is considered anomalous. *The specific criteria specified above were determined based on consideration of thermograms from 10 distinct SOA experiments. While these criteria should be robustly applicable to other FIGAERO-CIMS datasets, they can be adjusted depending on the specific application, data quality, and needs.*”

5. Similarly, lines 171-182 use a smoothing over 35 points. It is not clear how this smoothing is carried out, e.g. a mean of 35 points? Is it +17/-18, or -35, or +35 points that are smoothed? If the thermograms have peaks around 40 points wide are we seeing the mean of the value in the peak? If so how does this correlate with ‘retaining the peak shape’? I think this should be clarified.

A boxcar smoothing function is used, where the average is for the central point and the $(N-1)/2$ points before and after the point, where N is the number of points in the boxcar, here 35; this applies to all points that are at least $1+(N-1)/2$ points from the edges (the first and last point). We now state that a boxcar smoothing function is used. Points that are closer to the edges than $1+(N-1)/2$ points are excluded. Thus, the smoothed thermogram has $N-1$ fewer points than the unsmoothed, here a total of 34 out of 800. Since peaks are almost never observed at the start or end of the thermogram this exclusion of points does not impact the determination of the smoothed maximum for each thermogram. As to “retaining the peak

shape,” such smoothing only occasionally leads to lowering of the especially sharp peaks, however, will retain the location and shape of most peaks. We also show an example smoothed thermogram, compared to the original thermogram, below. This example is a reasonable sharp thermogram, to illustrate that the shape is nominally preserved, although the maximum is certainly reduced relative to the single-point maximum. We have not included such a figure in the manuscript, although have expanded the description of the smoothing in section 2.1.2.

“Normalization is achieved by dividing each point of the original thermogram by the thermogram maximum, where the maximum is determined after smoothing using a 35-point boxcar moving average with the end points excluded from the smoothed thermogram.”



6. Similarly, line 190 recommends a weighting of 4:1. How has this value been arrived at?

The value is empirical based on the information two different heating stages carry. We have also tried 1:1, 2:1 and 10:1 for the example data sets, yet 4:1 leads to the best clustering results. This value is ultimately user-defined depending on what kind of information they are trying to extract from the thermograms. We have expanded the discussion regarding this weighting factor.

“we recommend down-weighting the soaking period such that the ramping and soaking periods ultimately carry approximately 4:1 weight in the calculation of the ED. We have tested weighting of 1:1, 2:1 and 10:1. Weighting of 4:1 provides for the most robust clustering results for the example datasets.”

7. Lines 216-218 discuss the removal of noise. DBScan is considered deterministic for core points and noise. Noisy data would normally be identified and be members of clusters < minpts. Is there a danger that data identified with ‘high levels of noise’ is excluded when, despite the noise, it is similar enough to be included in a cluster?

The reviewer raises an important point about the treatment of noisy thermograms in the data processing stage. As we described in section 2.1.3, the noise level is inversely related to the signal level of an ion. Therefore, the excluded noisy data are usually unimportant and makes up only a small fraction of total signal. For all the example data sets, we have not discovered any excluded noisy thermogram carrying significant mass that also has a shape clearly similar to the clustered thermograms. We note that one reason for excluding the noisy thermograms is that this defines a criteria that helps guide inspection of the single-ion clusters (i.e., those having fewer than minpts). Similar to DBSCAN, the NSSC also has the ability to identify noisy data or outliers due to the insert of seed-sorting process. In most cases, NSSC with and without a pre-removal of noisy data gives identical clustering results. The user can choose to skip the treatment of noisy data, if desired. We have added the following as additional context regarding removal of noisy thermograms.

“This is especially the case for algorithms such as k-means and partitioning around medoids, which assign all the members to a cluster. Clustering methods that do not require assignment of all members, such as DBSCAN or our NSSC, are generally less sensitive to the influence of overly noisy members. However, we have found that the explicit exclusion of noisy thermograms up front serves to provide for more robust behavior and also removes the need to consider each noisy thermogram as a possible single-member cluster.”

8. Lines 226-235 appear to form the work being considered as ‘a novel variant of DBScan’. This describes DBScan with no alterations, except to force the order of data to consider ‘seed thermograms’ first. This is a pre-processing stage of the data, rather than a novel clustering algorithm. DBScan is deterministic if data order is preserved. If data order is not preserved then it is deterministic for core and noise, with only border points varying in some cases. I’d be interested to see how the border points vary to justify forcing the data order. Are the results generated by forcing the data order better simply because they [are] consistent with each run, whereas random initialisation is not? If so, is it possible to identify which thermograms change cluster and consider why?

As we discussed in the response to comment #2, the clustering algorithm of NSSC also differs from that of DBScan in the way they determine neighbors with a seed. We have expanded the discussion of the differences of NSSC and DBScan in section 2.3 and added the following statement to the description of the NSSC regarding the treatment of the seed: *“The seed does not evolve as neighbors are added to the cluster during this step.”* For the same dataset, NSSC generally results in more clusters. In NSSC there will not be border points, unlike DBScan, because there is only one time of scanning for each seed.

9. Lines 240-255 describe the DBScan algorithm. I am unclear how this varies from standard DBScan. Cycling through those thermograms identified as ‘seeds’ could equally be done by ordering the data by order of noise, then using the data, in order, to run standard DBScan.

As we discussed in the response to comment #2, the clustering algorithm of NSSC also differs from that of DBScan in the way they determine neighbors with a seed, besides the inclusion of seed-sorting and second-round of clustering processes.

10. Line 240 refers to Figure 2. This is not a suitable method for presenting an algorithm and a formal pseudo-code should be used. This may help clear up any confusion over the similarities or differences between DBScan and the method proposed.

It is our understanding that flowcharts are generally accepted as a way to present algorithms. Flowcharts for algorithm presentation have reasonably standardized symbols and structure, which we have endeavored to follow. When deciding between whether to include a flowchart or pseudocode, we considered whom we thought most likely to read and use this work and concluded the target audience is atmospheric chemists and other FIGAERO-CIMS users. As atmospheric chemists are generally familiar with flowcharts, but typically have less familiarity with pseudocode, we decided that presentation of our algorithm as a flowchart was preferable. We understand the reviewers point that clear distinction between DBScan and NSSC is necessary, and have added to the text to help in this regard, as discussed above. We note finally that the complete code of NSSC is also available at GitHub (<https://github.com/chrisappa/NSSC>).

11. Lines 256 – 269 describe a ‘second round’ of clustering. This generates new data for each cluster in the form of a ‘signal weighted average’, presumably of the cluster members. A thermogram that is within ϵ of the average, but not already clustered suggests that it is a border point, but below the *minpts* threshold for inclusion? (I am unclear how a thermogram can be within ϵ of the average, but not within ϵ of *minpts* of other cluster members?) This part appears to be a novel ‘second stage’, however I would not consider this to be a clustering algorithm in itself, but rather a post-processing step to tidy up ‘stragglers’, which is application specific.

As we discussed in the response to comment #2, the clustering algorithm of NSSC is different from that of DBScan in the way they determine neighbors with a seed. Therefore, there is no concept of “border points” in NSSC. Since the average thermogram of a cluster can be slightly different from seed thermogram, this second round of clustering is necessary to tidy up stragglers. The algorithm of the second round is more similar to how mean-shift method adjust the existing clusters. Therefore, we think this should be considered to be part of the algorithm.

12. Section 2.2.3 Describes a process for selecting an optimal ϵ value. The selection of ϵ is based on fuzzy terms such as ‘small’ and ‘near the maximum’. Figure 4 shows a clear value of ϵ in this case, is this the same for all analyses?

The reviewer raises an important point on the determination of optimal ϵ . In order to find the optimal ϵ , NSSC has to be run on the dataset for multiple times for a range of ϵ . We have shown in figure 4 as an example of how to determine the optimal ϵ . There are four parameters exhibiting different behaviors as a function of ϵ . At the current stage, we recommend users to determine optimal ϵ by visually comparing the values of these four parameters at different ϵ with the help of a figure such as figure 4. In the future, it would be ideal to find a way to determine the optimal ϵ automatically. We note that in leaving the decision about the optimal ϵ “fuzzy,” the approach here shares some similarity with other clustering algorithms for which, for example, the number of clusters or the minimum epsilon must be specified. It is our understanding that various approaches, such as “elbow plots” have been used to determine optimal parameters, but quite often these have an element of “fuzzy”-ness to them, as here. As to the reviewer’s question about whether there is a clear epsilon for all cases, we provided the determination plots for every experiment considered in the supplemental material. Fig. 4 is one example.

13. I am also unclear from section 2.2.3 whether this selection of optimal ϵ is generic to all future datasets of this type, or whether this optimal selection process is required for each new set of data?

The selection of ε is specific to each experiment and dataset. However, the guidance to determine an optimal value is generic. The determination plots, similar to Fig. 4, for each dataset are provided in the Supplemental Material.

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

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

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

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

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

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

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

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

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

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

15 ^c ToFwerk AG, Thun, Switzerland

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

17 Abstract

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

Deleted: novel

Deleted: and apply it

Deleted: NSSC

42 clustering will allow for more comprehensive understanding of the thermochemical properties
43 of OA.

44 **1. Introduction**

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

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

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

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

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

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

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

125 2. Clustering Method Description

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

Deleted: developed

Deleted: , named

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

132 **2.1. Data Preprocessing**

133 **2.1.1. Exclusion of anomalous thermograms**

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

Deleted: absolute

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

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

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

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

169 **2.1.2. Euclidean Distance**

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

176

$$177 \quad ED_{a,b} = \sum_n \sqrt{(a_n - b_n)^2} \quad (1)$$

178

179 An individual *ED* value is obtained for every pair of ions in the mass spectrum, resulting in an *n* x
180 *n* matrix of *ED* values with the diagonal elements all zero. The signal levels between individual
181 ions differ substantially, reflecting their relative abundances. Therefore, the *ED* calculation uses
182 normalized thermograms, allowing for comparison between thermogram profiles irrespective of
183 signal magnitude. Normalization is achieved by dividing each point of the original thermogram
184 by the thermogram maximum, [where the maximum is determined](#) after smoothing using a
185 35-point [boxcar](#) moving average [with the end points excluded from the smoothed thermogram](#).
186 Use of the smoothed maximum instead of the unsmoothed maximum reduces the influence of

187 noise on normalization. In the FIGAERO datasets used in this study, a typical thermogram has a
188 temperature resolution of $\Delta T \sim 0.7$ °C during the ramping period, and a 35-point smooth
189 corresponds to smoothing over ~ 24.5 °C. Typical FIGAERO thermograms exhibit peaks ca. 40 °C
190 wide, and thus a 35-point smoothing retains the main peak shape while reducing the influence
191 of noise. In the constant temperature part of the thermogram (soaking period), signal levels
192 change slowly with time, on average less than 5 % for a 35 points (~ 3 minutes) period, so a
193 35-point smoothing is also appropriate. We note that the unsmoothed profiles are those that are
194 normalized; smoothing relates only to determining the maximum signal values used for
195 normalization.

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

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

214 **2.1.3. Dealing with noise**

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

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

229 Discussed further in section 2.3, clustering algorithms often perform poorly when overly
230 noisy data are included in the clustering. This is especially the case for algorithms such as k-means
231 and partitioning around medoids, which assign all the members to a cluster. [Clustering methods](#)
232 [that do not require assignment of all members, such as DBSCAN or our NSSC, are generally less](#)
233 [sensitive to the influence of overly noisy members. However, we have found that the explicit](#)
234 [exclusion of noisy thermograms up front serves to provide for more robust behavior and also](#)
235 [removes the need to consider each noisy thermogram as a possible single-member cluster.](#) The
236 inclusion of overly noisy peaks might obscure the underlying structure of clustered thermograms.
237 Noisy thermograms are identified as follows. First, the 5% of ions having the lowest noise are
238 identified. The ξ value of the noisiest ion from this subset of low-noise ions is defined as the
239 reference noise level, ξ_{ref} . Small differences in the choice of this threshold (e.g. using the lowest
240 7% of ions) do not materially influence the results. Ions for which $\xi_n > 3 \cdot \xi_{ref}$ are considered noisy
241 and excluded from the initial clustering. For the experiments we examined, there are 88-120 out
242 of ~ 300 ions left after noise screening, contributing 83.5% - 92.5% to the total particle mass.

243 2.2. Noise-sorted Scanning Clustering (NSSC)

244 2.2.1. Algorithm description

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

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

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

Deleted: ,

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

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

291 **2.2.2. Post-clustering Processes**

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

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

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

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

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

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

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

334 2.2.3. Choosing the optimal ϵ

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

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

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

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

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

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

369

$$370 \quad R_{interClst} = \frac{ED_{seed,avg}}{\varepsilon} = \frac{\sum_{i=1}^{N_{c,total}} \sum_{j=1}^{N_{c,total}} ED_{seed,i,j}}{N_{c,total}(N_{c,total}-1)\varepsilon} \quad (2)$$

371

372 and $ED_{seed,i,j}$ is the distance between the seeds for the different clusters i and j and $N_{c,total} = N_c +$
 373 $N_{c,one}$. For a 2D data set, the seed can be visualized as the center of a circle and ε the radius of
 374 the circle. Thus, when $ED_{seed,i,j}/\varepsilon < 2$, the two circles defining the boundaries of these two clusters
 375 have overlapping areas. Good separation (i.e. cluster dissimilarity) is indicated when $ED_{seed,i,j}/\varepsilon >$
 376 2 . Although our data set is more than two dimensions, this illustrates the idea of establishing the
 377 level of similarity (or dissimilarity) between clusters, i.e., the extent to which they are unique. We
 378 recommend selecting an ε that results in $R_{interClst} \geq 2$, when possible.

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

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

388 **2.3. Alternative Clustering Methods**

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

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

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

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

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

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

445 **3. FIGAERO Measurements and Experiments**

446 **3.1. Instrument and experiment description**

447 The FIGAERO-CIMS instrument has been described previously in detail (Lee et al., 2014;
448 Lopez-Hilfiker et al., 2014). A brief description is provided here, with some additional details in
449 the Supplemental Material. The FIGAERO-CIMS measures the evolved gases from filter-collected
450 particles during temperature programmed thermal desorption. Thermal desorption of particles
451 occurs in two-stages: a “ramping” and “soaking” period. During ramping, the temperature
452 increases from room temperature to 200 °C, typically at 10 °C min⁻¹. Most OA mass desorbs
453 during the ramping stage. The temperature is held at 200 °C for ca. 30–40 mins during the soaking
454 period to facilitate evaporation of the remaining, low-volatility organic mass from the filter. The
455 evolved gas-phase compounds are measured using CIMS with the iodide (I⁻) reagent ion,
456 appropriate for characterization of generally highly oxygenated components comprising most
457 secondary organic aerosol (Lopez-Hilfiker et al., 2016; Isaacman-VanWertz et al., 2017; Lee et al.,
458 2018). The resulting signal or mass concentration versus temperature (or equivalently time)
459 curves for each ion constitute a thermogram. All individual thermograms are background
460 corrected by subtracting the observed thermograms from appropriate blank experiments. The
461 overall bulk thermogram is obtained by summing together the individual thermograms.

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

471 3.2. General interpretation of FIGAERO-CIMS thermograms

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

4. Example Applications

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

4.1. α -pinene + OH SOA

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

527 A total of 11 clusters are identified with no one-member clusters. The unweighted and
528 mass-weighted average thermograms for each cluster are shown along with the thermograms of
529 individual members in **Figure 5a**. The differences between weighted and unweighted average
530 clusters are negligible, in general. Clusters are organized and numbered (as Clst#*N*) from low to
531 high T_{m50} , with deeper to shallower downslope. Clst#1 through Clst#6 all have a clear peak below
532 120 °C, but with different peak widths and downslopes. Clst#7 and Clst#8 are a bit noisier with
533 only a few members each, exhibiting a sharp upslope and shallow downslope. Clst#9 has a very
534 broad peak. Clst#10 peaks at around 150 °C after an initial rise and temporary plateau. Clst#11
535 exhibits behavior somewhat like Clst#10, but with a peak that occurs just into the soaking period,
536 evident if viewed in time space, at 200 °C with a rapid drop afterwards.

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

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

555 Interpretation of previous FIGAERO-CIMS studies have largely focused on the behavior of
556 the bulk thermogram or of several major ions or sums of ions based on common factors such as
557 the number of carbon atoms (Lopez-Hilfiker et al., 2016; D'Ambro et al., 2017; D'Ambro et al.,
558 2018; Stolzenburg et al., 2018; Wang and Ruiz, 2018; Joo et al., 2019). The normalized
559 thermograms of the top five ions contributing most to the total mass for the experiments here
560 are shown in **Figure 5c**, along with the bulk thermogram. Together these five ions make up nearly
561 30% of the total mass, and exhibit very similar thermogram shapes to each other and to the bulk
562 thermogram and belong solely to either Clst#2 or Clst#3. Thus, examining these ions only would
563 capture only a fraction of the overall diversity in thermal behaviors. The clustering method
564 developed here provides a means to investigate more comprehensively the variability in volatility
565 between aerosol components.

566 **4.2. Δ -3-carene + OH SOA**

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

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

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

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

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

597 There are ultimately three major differences between the two SOA systems. For one, there
598 is a different relationship between fractional contribution and cluster number (and thus $T_{m,50}$)
599 between the two. Secondly, the α -pinene SOA contains ions with especially narrow peaks at ca.
600 100 °C (i.e., Clst#7 & 8), that are not observed with Δ -3-carene SOA (compare **Figure 5** with **Figure**
601 **7**). Lastly, the thermograms of the top five ions for Δ -3-carene SOA differ to a greater extent than
602 for α -pinene SOA. Although we are unable to determine the reasons for these differences here,
603 this illustrates the potential for clustering to help identify and understand differences between
604 different SOA systems.

605 **4.3. α -pinene + OH + NO SOA**

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

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

618 **4.3.1. Single Clustering**

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

628 Mass spectra for the three experiments are compared in **Figure 8a** and the bulk
629 thermograms shown in **Figure 8b** and c. The 5 ppb NO and 10 ppb NO SOA mass spectra are
630 nearly identical. The mass spectrum for the 25 ppb NO experiment, however, exhibits a notable
631 shift of the most abundant ions towards lower m/z . The bulk thermograms for the 5 ppb and 10
632 ppb NO experiments are nearly identical, peaking near 80 °C. The 25 ppb NO bulk thermogram
633 similarly peaks near 80 °C, but exhibits a much slower decay as temperature increases further.
634 Additionally, the change in slope at the transition from the ramping to soaking period is more
635 pronounced in the 25 ppb NO experiment. Overall, a greater fraction of the mass desorbs above
636 100 °C and during the soaking period for the 25 ppb NO experiment compared to lower-NO
637 experiments.

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

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

659 **4.3.2. Multiple Clustering**

660 With multiple clustering, each experiment was processed and clustered independently,
661 with experiment-specific ξ_{ref} , N_c , and ε , among other parameters (**Figure S4** and **Table 3**). The
662 clustered thermograms from the three experiments are compared in **Figure 10a-c**. The number
663 of clusters identified increases with NO concentration. Comparison between the shapes of the
664 clusters from the 5 ppb NO (**Figure 10a**) and 10 ppb NO (**Figure 10b**) experiments indicates
665 generally similar types of thermograms, consistent with the single clustering method. Ten of the
666 11 total 10 ppb clusters match with a 5 ppb cluster. The one additional, unique cluster at 10 ppb

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

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

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

687 The percent contribution of different clusters to total mass, along with the noise-filtered
688 and unclustered ions, differ between experiments (**Figure 10d**). Note that for the multiple
689 clustering method, clusters having the same index number are not necessarily directly
690 comparable between experiments because different sets of ions are included. For example, while
691 Clst#1 in the 5 ppb and 10 ppb NO experiments are comparable, the most similar cluster in the
692 25 ppb experiment is Clst#2. Nonetheless, there are some common features shared by the same,
693 or closely indexed, clusters. For example, Clst#1 – 4 in all three experiments exhibit a narrow,
694 single peak with the peak temperature below 120 °C. The mass contribution of Clst#1 – 4 is similar
695 between the 5 and 10 ppb NO experiment, but ~15% lower in the 25 ppb NO experiment. Clusters

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

703 **4.4. α -pinene + O₃ SOA**

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

711 **4.4.1. Single Clustering**

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

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

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

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

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

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

765 The extent of mass loss with isothermal evaporation differs between clusters. In general,
766 clusters that exhibit larger changes in shape have greater total mass loss, although with variability
767 (**Figure S6c**). Consequently, the mass contributions of the clusters evolve with isothermal
768 evaporation (**Figure 12b**). The contribution of Clst#1 decreases significantly and most notably as
769 wait time increases. The most prominent ion in the no-wait experiment, $C_{10}H_{14}O_6$, is grouped in
770 Clst#1. The continuous mass loss of Clst#1 indicates the rapid evaporation of its members. The
771 mass contributions of the other clusters that exhibited similar changes in shape as Clst#1 (Clst#3,
772 5, and 6) remain comparably constant, although with Clst#3 decreasing slightly. The relative
773 abundances of the clusters for which the thermograms shapes changed negligibly (Clst#7 and 11)
774 increase continually, implying of the slowest evaporation of the ions in these two clusters in the
775 24-hr evaporation period.

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

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

787 **4.4.2. Multiple Clustering**

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

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

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

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

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

827 **5. Conclusions**

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

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

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

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

Deleted: properties of organic aerosol in further detail.

Formatted: Font color: Accent 1

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

872 **6. Data Availability**

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

877 **7. Author Contributions**

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

881 **8. Acknowledgements**

882 This work was supported by the National Science Foundation under Grant No. ATM-
883 1151062. The experimental work described here was supported by the U.S. Department of
884 Energy ASR grants DE-SC0011791 and DE-SC0018221. E.L.D. was supported by the National
885 Science Foundation Graduate Research Fellowship (grant no. DGE-1256082) and S.S. was
886 supported by the Academy of Finland (grant nos. 272041 and 310682). The SOAFFEE campaign
887 was done at Pacific Northwest National Laboratory, supported by the U.S. Department of Energy
888 (DOE) Office of Science, Office of Biological and Environmental Research, as part of the
889 Atmospheric Systems Research (ASR) program. PNNL is operated for DOE by Battelle Memorial
890 Institute under contract DE-AC05-76RL01830.

891 9. References

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

933 Spectrometer, *Atmos Environ*, 61, 316-326, <https://doi.org/10.1016/j.atmosenv.2012.07.054>,
934 2012.

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

937 Gonzalez, T. F.: Clustering to Minimize the Maximum Intercluster Distance, *Theor Comput Sci*, 38,
938 293-306, [https://doi.org/10.1016/0304-3975\(85\)90224-5](https://doi.org/10.1016/0304-3975(85)90224-5), 1985.

939 Hamilton, J. F., Webb, P. J., Lewis, A. C., Hopkins, J. R., Smith, S., and Davy, P.: Partially oxidised
940 organic components in urban aerosol using GCXGC-TOF/MS, *Atmospheric Chemistry and Physics*,
941 4, 1279-1290, <https://doi.org/10.5194/acp-4-1279-2004>, 2004.

942 Huang, W., Saathoff, H., Pajunoja, A., Shen, X. L., Naumann, K. H., Wagner, R., Virtanen, A., Leisner,
943 T., and Mohr, C.: alpha-Pinene secondary organic aerosol at low temperature: chemical
944 composition and implications for particle viscosity, *Atmospheric Chemistry and Physics*, 18, 2883-
945 2898, <https://doi.org/10.5194/acp-18-2883-2018>, 2018.

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

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

954 Kirchner, U., Vogt, R., Natzeck, C., and Goschnick, J.: Single particle MS, SNMS, SIMS, XPS, and
955 FTIR spectroscopic analysis of soot particles during the AIDA campaign, *Journal of Aerosol Science*,
956 34, 1323-1346, [https://doi.org/10.1016/S0021-8502\(03\)00362-8](https://doi.org/10.1016/S0021-8502(03)00362-8), 2003.

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

961 Lee, A. K. Y., Willis, M. D., Healy, R. M., Onasch, T. B., and Abbatt, J. P. D.: Mixing state of
962 carbonaceous aerosol in an urban environment: single particle characterization using the soot
963 particle aerosol mass spectrometer (SP-AMS), *Atmospheric Chemistry and Physics*, 15, 1823-
964 1841, <https://doi.org/10.5194/acp-15-1823-2015>, 2015.

965 Lee, B., Lopez-Hilfiker, F. D., D'Ambro, E. L., Zhou, P. T., Boy, M., Petaja, T., Hao, L. Q., Virtanen,
966 A., and Thornton, J. A.: Semi-volatile and highly oxygenated gaseous and particulate organic
967 compounds observed above a boreal forest canopy, *Atmospheric Chemistry and Physics*, 18,
968 11547-11562, <https://doi.org/10.5194/acp-18-11547-2018>, 2018.

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

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

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

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

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

987 Liu, S., Shilling, J. E., Song, C., Hiranuma, N., Zaveri, R. A., and Russell, L. M.: Hydrolysis of
988 Organonitrate Functional Groups in Aerosol Particles, *Aerosol Science and Technology*, 46, 1359-
989 1369, <https://doi.org/10.1080/02786826.2012.716175>, 2012.

990 Liu, S., Russell, L. M., Sueper, D. T., and Onasch, T. B.: Organic particle types by single-particle
991 measurements using a time-of-flight aerosol mass spectrometer coupled with a light scattering
992 module, *Atmospheric Measurement Techniques*, 6, 187-197, [https://doi.org/10.5194/amt-6-](https://doi.org/10.5194/amt-6-187-2013)
993 [187-2013](https://doi.org/10.5194/amt-6-187-2013), 2013.

994 Lopez-Hilfiker, F. D., Mohr, C., Ehn, M., Rubach, F., Kleist, E., Wildt, J., Mentel, T. F., Lutz, A.,
995 Hallquist, M., Worsnop, D., and Thornton, J. A.: A novel method for online analysis of gas and
996 particle composition: description and evaluation of a Filter Inlet for Gases and AEROSols
997 (FIGAERO), *Atmospheric Measurement Techniques*, 7, 983-1001, [https://doi.org/10.5194/amt-](https://doi.org/10.5194/amt-7-983-2014)
998 [7-983-2014](https://doi.org/10.5194/amt-7-983-2014), 2014.

999 Lopez-Hilfiker, F. D., Mohr, C., Ehn, M., Rubach, F., Kleist, E., Wildt, J., Mentel, T. F., Carrasquillo,
1000 A. J., Daumit, K. E., Hunter, J. F., Kroll, J. H., Worsnop, D. R., and Thornton, J. A.: Phase partitioning
1001 and volatility of secondary organic aerosol components formed from α -pinene ozonolysis and OH
1002 oxidation: the importance of accretion products and other low volatility compounds,
1003 *Atmospheric Chemistry and Physics*, 15, 7765-7776, <https://doi.org/10.5194/acp-15-7765-2015>,
1004 2015.

1005 Lopez-Hilfiker, F. D., Mohr, C., D'Ambro, E. L., Lutz, A., Riedel, T. P., Gaston, C. J., Iyer, S., Zhang,
1006 Z., Gold, A., Surratt, J. D., Lee, B. H., Kurten, T., Hu, W. W., Jimenez, J., Hallquist, M., and Thornton,
1007 J. A.: Molecular Composition and Volatility of Organic Aerosol in the Southeastern U.S.:
1008 Implications for IEPOX Derived SOA, *Environ Sci Technol*, 50, 2200-2209,
1009 <https://doi.org/10.1021/acs.est.5b04769>, 2016.

1010 Mohr, C., Lopez-Hilfiker, F. D., Yli-Juuti, T., Heitto, A., Lutz, A., Hallquist, M., D'Ambro, E. L.,
1011 Rissanen, M. P., Hao, L. Q., Schobesberger, S., Kulmala, M., Mauldin, R. L., Makkonen, U., Sipila,
1012 M., Petaja, T., and Thornton, J. A.: Ambient observations of dimers from terpene oxidation in the
1013 gas phase: Implications for new particle formation and growth, *Geophysical Research Letters*, 44,
1014 2958-2966, <https://doi.org/10.1002/2017gl072718>, 2017.

1015 Murphy, D. M., Middlebrook, A. M., and Warshawsky, M.: Cluster analysis of data from the
1016 Particle Analysis by Laser Mass Spectrometry (PALMS) instrument, *Aerosol Science and*
1017 *Technology*, 37, 382-391, <https://doi.org/10.1080/02786820300971>, 2003.

1018 Pinero-Garcia, F., Ferro-Garcia, M. A., Chham, E., Cobos-Diaz, M., and Gonzalez-Rodelas, P.: A
1019 cluster analysis of back trajectories to study the behaviour of radioactive aerosols in the south-

1020 east of Spain, *J Environ Radioactiv*, 147, 142-152, <https://doi.org/10.1016/j.jenvrad.2015.05.029>,
1021 2015.

1022 Praske, E., Otkjaer, R. V., Crouse, J. D., Hethcox, J. C., Stoltz, B. M., Kjaergaard, H. G., and
1023 Wennberg, P. O.: Atmospheric autoxidation is increasingly important in urban and suburban
1024 North America, *P Natl Acad Sci USA*, 115, 64-69, <https://doi.org/10.1073/pnas.1715540115>, 2018.

1025 Rebotier, T. P., and Prather, K. A.: Aerosol time-of-flight mass spectrometry data analysis: A
1026 benchmark of clustering algorithms, *Anal Chim Acta*, 585, 38-54,
1027 <https://doi.org/10.1016/j.aca.2006.12.009>, 2007.

1028 Reitz, P., Zorn, S. R., Trimborn, S. H., and Trimborn, A. M.: A new, powerful technique to analyze
1029 single particle aerosol mass spectra using a combination of OPTICS and the fuzzy c-means
1030 algorithm, *Journal of Aerosol Science*, 98, 1-14, <https://doi.org/10.1016/j.jaerosci.2016.04.003>,
1031 2016.

1032 Roth, A., Schneider, J., Klimach, T., Mertes, S., van Pinxteren, D., Herrmann, H., and Borrmann, S.:
1033 Aerosol properties, source identification, and cloud processing in orographic clouds measured by
1034 single particle mass spectrometry on a central European mountain site during HCCT-2010,
1035 *Atmospheric Chemistry and Physics*, 16, 505-524, <https://doi.org/10.5194/acp-16-505-2016>,
1036 2016.

1037 Schobesberger, S., D'Ambro, E. L., Lopez-Hilfiker, F. D., Mohr, C., and Thornton, J. A.: A model
1038 framework to retrieve thermodynamic and kinetic properties of organic aerosol from
1039 composition-resolved thermal desorption measurements, *Atmospheric Chemistry and Physics*,
1040 18, 14757-14785, <https://doi.org/10.5194/acp-18-14757-2018>, 2018.

1041 Song, X. H., Hopke, P. K., Fergenson, D. P., and Prather, K. A.: Classification of single particles
1042 analyzed by ATOFMS using an artificial neural network, *ART-2A, Anal Chem*, 71, 860-865,
1043 <https://doi.org/10.1021/ac9809682>, 1999.

1044 Stolzenburg, D., Fischer, L., Vogel, A. L., Heinritzi, M., Schervish, M., Simon, M., Wagner, A. C.,
1045 Dada, L., Ahonen, L. R., Amorim, A., Baccarini, A., Bauer, P. S., Baumgartner, B., Bergen, A., Bianchi,
1046 F., Breitenlechner, M., Brilke, S., Mazon, S. B., Chen, D. X., Dias, A., Draper, D. C., Duplissy, J.,
1047 Haddad, I., Finkenzeller, H., Frege, C., Fuchs, C., Garmash, O., Gordon, H., He, X., Helm, J.,
1048 Hofbauer, V., Hoyle, C. R., Kim, C., Kirkby, J., Kontkanen, J., Kuerten, A., Lampilahti, J., Lawler, M.,
1049 Lehtipalo, K., Leiminger, M., Mai, H., Mathot, S., Mentler, B., Molteni, U., Nie, W., Nieminen, T.,
1050 Nowak, J. B., Ojdanic, A., Onnela, A., Passananti, M., Petaja, T., Quelever, L. L. J., Rissanen, M. P.,
1051 Sarnela, N., Schallhart, S., Tauber, C., Tome, A., Wagner, R., Wang, M., Weitz, L., Wimmer, D.,
1052 Xiao, M., Yan, C., Ye, P., Zha, Q., Baltensperger, U., Curtius, J., Dommen, J., Flagan, R. C., Kulmala,
1053 M., Smith, J. N., Worsnop, D. R., Hansel, A., Donahue, N. M., and Winkler, P. M.: Rapid growth of
1054 organic aerosol nanoparticles over a wide tropospheric temperature range, *P Natl Acad Sci USA*,
1055 115, 9122-9127, <https://doi.org/10.1073/pnas.1807604115>, 2018.

1056 Takahama, S., Gilardoni, S., Russell, L. M., and Kilcoyne, A. L. D.: Classification of multiple types
1057 of organic carbon composition in atmospheric particles by scanning transmission X-ray
1058 microscopy analysis, *Atmos Environ*, 41, 9435-9451,
1059 <https://doi.org/10.1016/j.atmosenv.2007.08.051>, 2007.

1060 Wang, D. S., and Ruiz, L. H.: Chlorine-initiated oxidation of n-alkanes under high-NO_x conditions:
1061 insights into secondary organic aerosol composition and volatility using a FIGAERO-CIMS,
1062 *Atmospheric Chemistry and Physics*, 18, 15535-15553, [https://doi.org/10.5194/acp-18-15535-
1063 2018](https://doi.org/10.5194/acp-18-15535-2018), 2018.

1064 Wegner, T., Hussein, T., Hameri, K., Vesala, T., Kulmala, M., and Weber, S.: Properties of aerosol
1065 signature size distributions in the urban environment as derived by cluster analysis, Atmos
1066 Environ, 61, 350-360, <https://doi.org/10.1016/j.atmosenv.2012.07.048>, 2012.

1067 Zhao, W. X., Hopke, P. K., and Prather, K. A.: Comparison of two cluster analysis methods using
1068 single particle mass spectra, Atmos Environ, 42, 881-892,
1069 <https://doi.org/10.1016/j.atmosenv.2007.10.024>, 2008.

1070 Zhao, Y., Thornton, J. A., and Pye, H. O. T.: Quantitative constraints on autoxidation and dimer
1071 formation from direct probing of monoterpene-derived peroxy radical chemistry, P Natl Acad Sci
1072 USA, 115, 12142-12147, <https://doi.org/10.1073/pnas.1812147115>, 2018.

1073 Zhou, L. M., Hopke, P. K., and Venkatachari, P.: Cluster analysis of single particle mass spectra
1074 measured at Flushing, NY, Anal Chim Acta, 555, 47-56, <https://doi.org/10.1016/j.aca.2005.08.061>,
1075 2006.
1076

1077

1078 **10. Tables**1079 **Table 1.** Details of SOA formation and chamber conditions for all the example SOA systems.

Exp #	Precursor		Oxidant		Seeds		UV	T (°C)	RH (%)	NO ^{#5} (ppb)	M _p ^{#&} (µg/m ³)	FIGAERO Operation _s
	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	α-pinene	10	OH (H ₂ O ₂)	1.0	AS	50	On	25	50	5	8.3	Normal
3b										10	9.2	
3c										25	9.1	
4a	α-pinene	10	O ₃	0.1	PS ^{&&}	50	Off	25	80	-	4.0	Normal
4b												1 h wait
4c												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

#⁵ 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

s 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

1080

1081

1082 **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

1083
1084

1085 **Table 3.** Parameters and thresholds used for the data processing and noise-sorted scanning clustering for
 1086 all the example experiments.

Expt #	SOA type	Pre-processing						Clustering					
		N_{total}	$N_{anomalous}$	$N_{filtered}$	$f_{m,filtered}$	ξ_{ref}	$f_{m,ref}$	ε	N_c	$N_{c,one}$	$f_{m,unclustered}$	$R_{interClist}$	
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	298	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	298	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			10	185	11.5	0.025	0.42	2.2	10	2	0.67	2.28	
4b	Single	312	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	312	1	191	11.4	0.025	0.41	2.2	11	2	1.04
4b	0	210			16.5	0.044	0.41	3.3	8	4	0.00	2.02	
4c	5	205			14.3	0.048	0.42	3.1	9	2	1.06	1.66	
4d	3	203			12.8	0.055	0.39	3.3	8	1	2.50	1.80	
4e	3	213			16.1	0.053	0.41	3.4	7	2	0.98	1.97	

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

$N_{anomalous}$ – Number of anomalous ions

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

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

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

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

ε – distance criterion

N_c – Number of clusters determined with two or more members

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

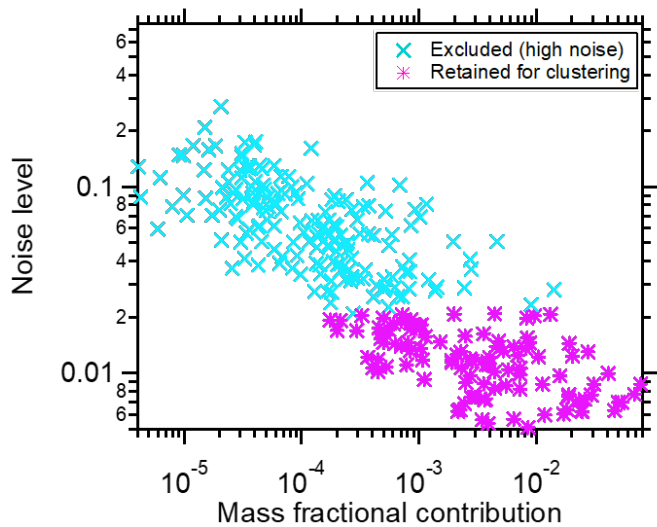
$f_{m,unclustered}$ – Mass fraction of unclustered ions, expressed in %

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

1087

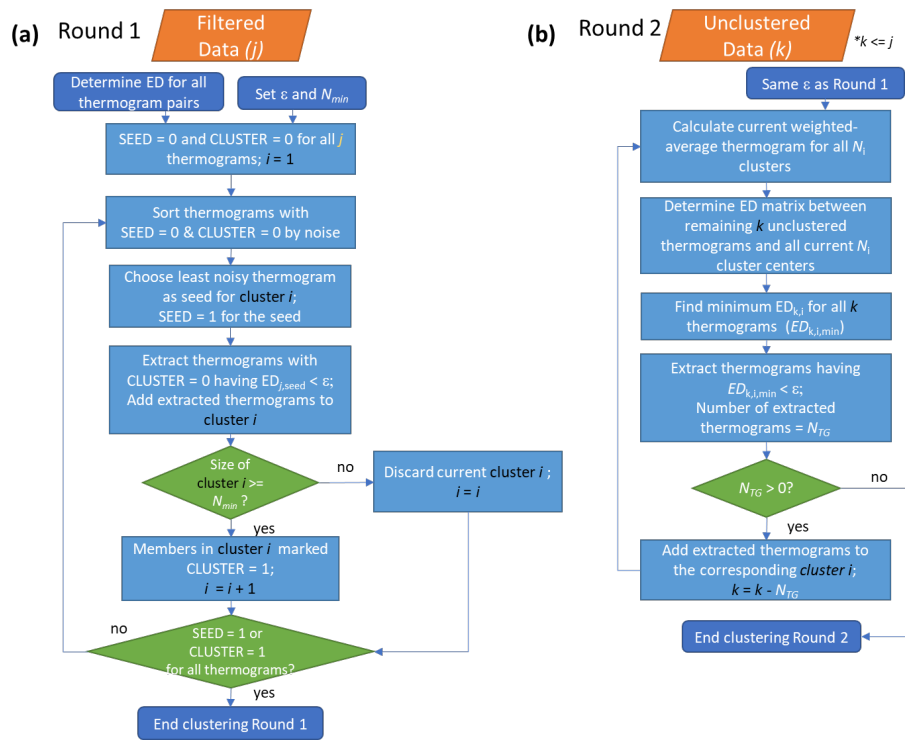
1088

1089 **11. Figures**



1090 **Figure 1:** The relationship between thermogram noise levels and the fractional contributions of the
1091 corresponding ions to total mass, for α -pinene + OH SOA. The noise threshold, $\xi_{ref} = 0.021$ and is used to
1092 distinguish high-noise thermograms (cyan markers) from thermograms having acceptable noise levels
1093 (pink markers).
1094

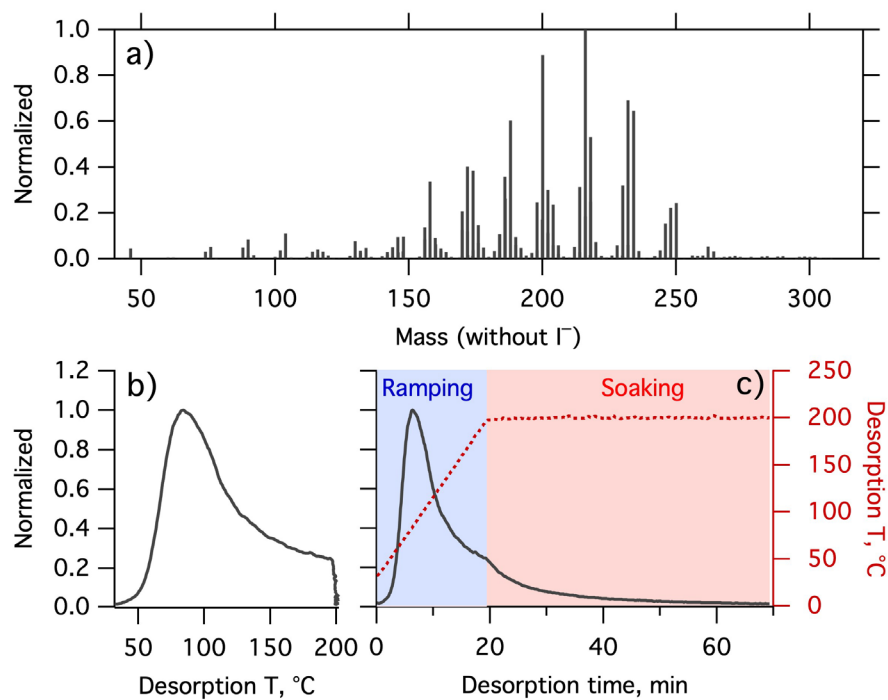
1095



1096
 1097 **Figure 2:** Flow of the noise-sorted scanning clustering. There are two rounds of clustering. (a) Round 1:
 1098 The ED between all thermogram pairs are calculated and two parameters, ϵ and N_{min} , are set. Each
 1099 thermogram is initialized with state $SEED = 0$ and $CLUSTER = 0$. Only thermograms with $SEED = 0$ and
 1100 $CLUSTER = 0$ can serve as seeds, while thermograms with $CLUSTER = 0$ can be added to new clusters. The
 1101 procedure terminates when all the thermograms are marked either $SEED = 1$ or $CLUSTER = 1$. (b) Round
 1102 2: Seeds are specified as the weighted-average thermogram for each cluster, and any remaining
 1103 unclustered thermograms from Round 1 are potentially added to these clusters. With the indexing, j refers
 1104 to the total number of thermograms, i to the number of clusters, and k to the number of unclustered
 1105 thermograms after Round 1.

1106

1107

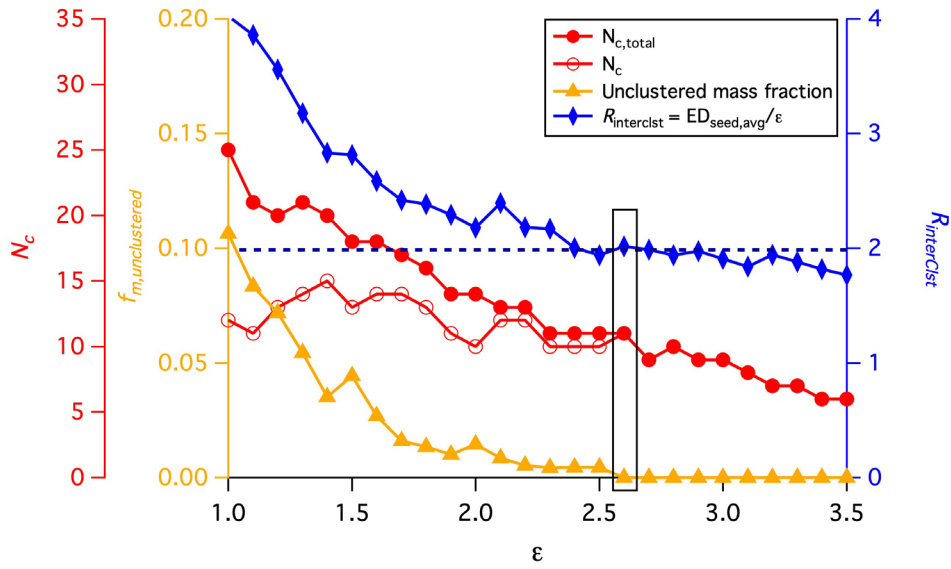


1108
1109

1110 **Figure 3.** (a) Mass spectrum of α -pinene + OH SOA measured by FIGAERO-CIMS. The mass excludes iodine.
1111 (b) Normalized thermogram of the bulk SOA versus temperature. (c) Normalized thermogram of the bulk
1112 SOA versus time (black line) and the variation in desorption temperature with time (dark red dashed line).
1113 The long tail during the soaking period is evident when the thermogram is considered in time space. The
1114 light blue shaded area denotes the ramping period and the pink shaded area the soaking period.

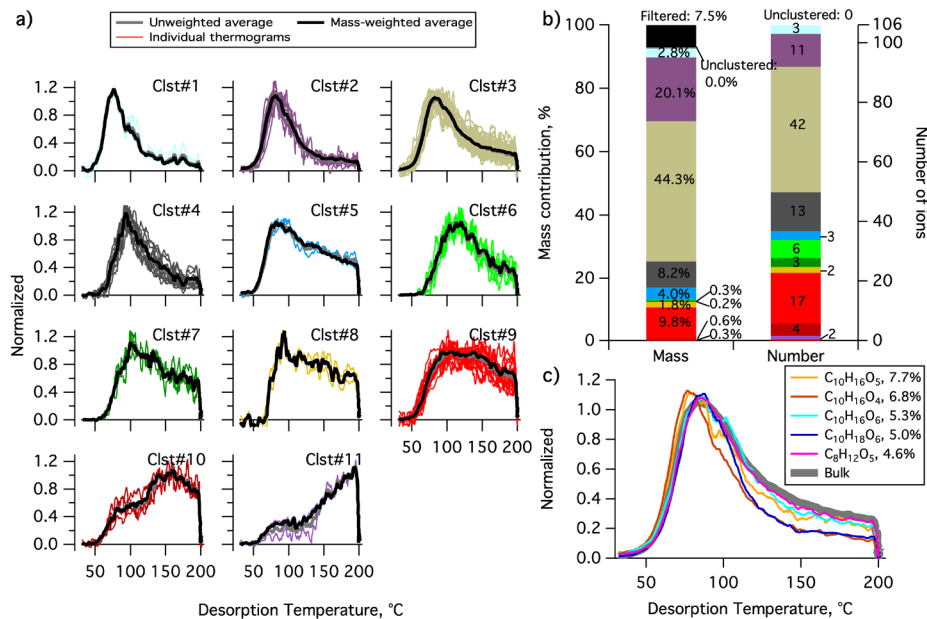
1115

1116



1117 **Figure 4.** The variation of four parameters, N_c , $N_{c,total}$, $f_{m,unclustered}$ and $R_{interClst}$ as a function of the distance
1118 criterion ϵ . The black horizontal dashed line guides the judgement for $R_{interClst} \geq 2$. The values highlighted
1119 by a rectangle are the values corresponding to the optimal ϵ used for the clustering analysis.
1120

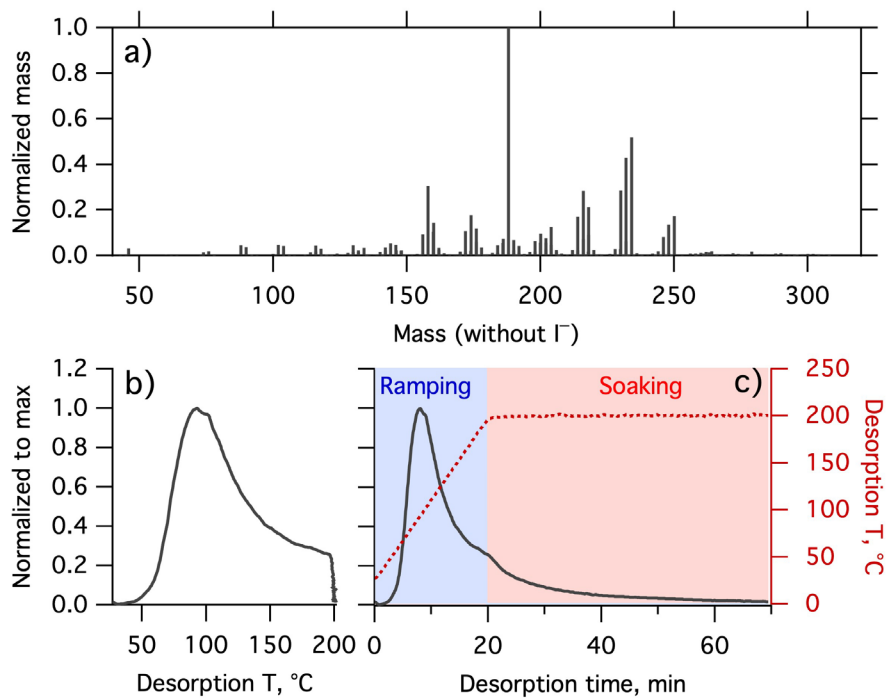
1121



1122
 1123 **Figure 5.** Clustering results for α -pinene + OH SOA. (a) Unweighted average thermograms (bold grey lines),
 1124 mass-weighted average thermograms (bold black lines) and individual members (colored lines) of the 11
 1125 clusters identified. (b) Percentage contribution of each cluster to the total mass, as well as the filtered out
 1126 and unclustered mass percentage (left bar), and the number of ions in each cluster and the unclustered
 1127 number of ions (right bar). (c) Thermograms of the top 5 ions in terms of mass contribution. The cluster
 1128 colors are consistent between (a) and (b).

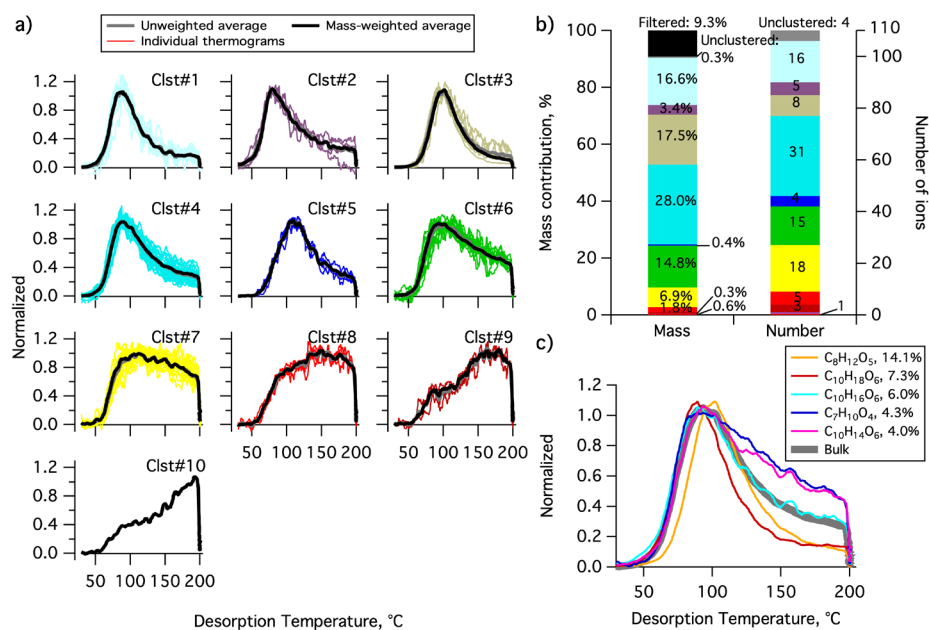
1129

1130



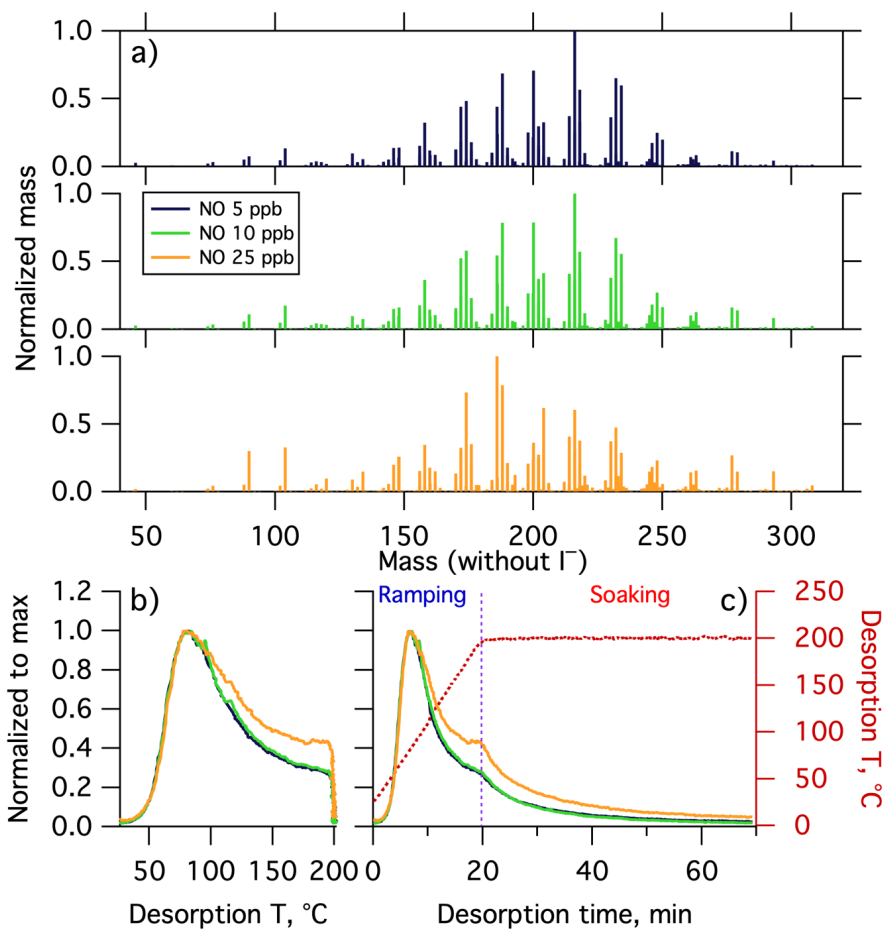
1131
1132 **Figure 6.** Same as Figure 3, but for Δ -3-carene + OH SOA. (a) SOA mass spectrum measured by
1133 FIGAERO-CIMS. The mass excludes iodine. The normalized thermogram of the bulk SOA versus (b)
1134 temperature and (c) time. In (c) the light blue shaded area denotes the ramping period and the pink
1135 shaded area the soaking period.

1136



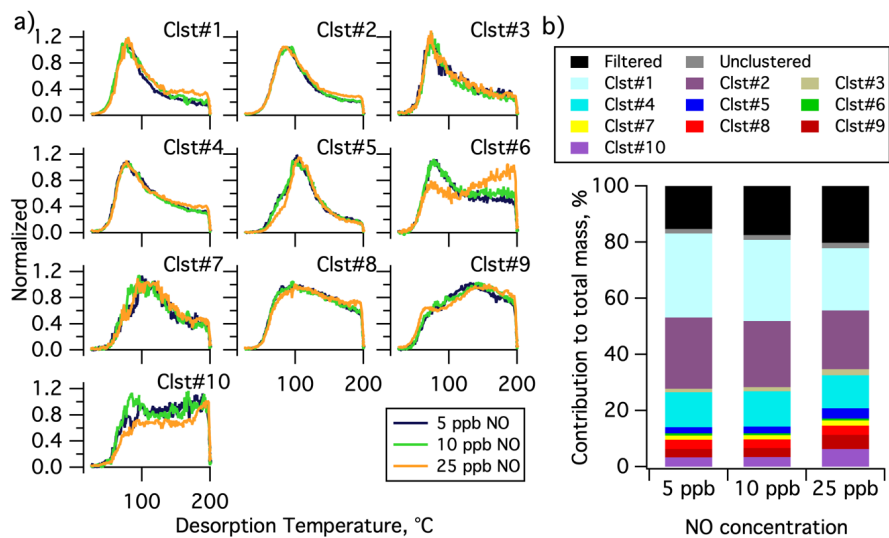
1137
 1138 **Figure 7.** Same as Figure 5, but for Δ -3-carene + OH SOA. (a) Unweighted average thermograms (bold grey
 1139 lines), mass-weighted average thermograms (bold black lines) and individual members (colored lines) of
 1140 the ten clusters identified. (b) Percentage contribution of each cluster to the total mass, as well as the
 1141 filtered out and unclustered mass percentage (left bar) and number of ions in each cluster and the
 1142 unclustered number of ions (right bar). (c) Thermograms of the top 5 ions in terms of mass contribution.
 1143 The cluster colors are consistent between (a) and (b).

1144



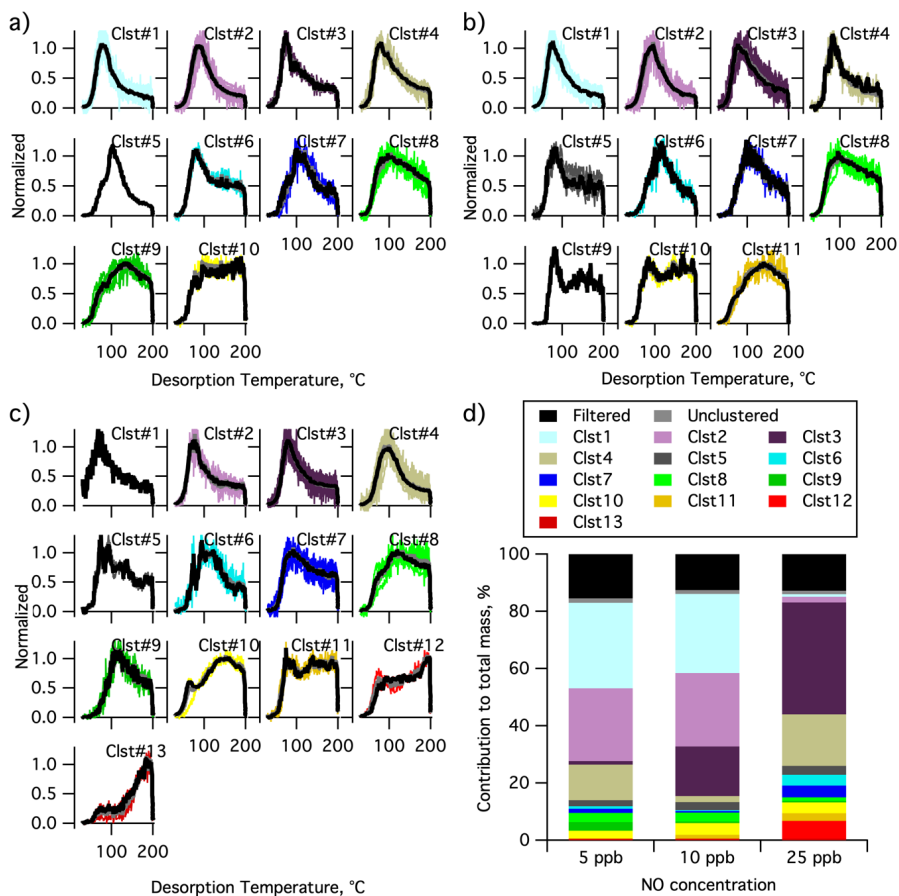
1145
 1146 **Figure 8.** (a) Mass spectra of α -pinene + OH SOA formed with different NO concentrations, normalized to
 1147 the most abundant ions mass concentration. The mass excludes iodine. Normalized thermograms of the
 1148 bulk SOA versus (b) temperature and (c) desorption time, with the desorption temperature shown in dark
 1149 red dashed line. The vertical purple dashed line delineates between ramping and soaking. In all the panels,
 1150 colors correspond to the NO concentration (see legend).

1151



1152
 1153 **Figure 9.** Single clustering results for α -pinene + OH SOA as a function of NO concentration. (a)
 1154 Comparison of the normalized, weighted average thermograms of the ten clusters for the 5 ppb NO (navy),
 1155 10 ppb NO (green) and 25 ppb NO (orange) experiments. (b) Contribution of each cluster to the total mass,
 1156 including the contribution from filtered out ions (black) and unclustered ions (gray). The total mass is
 1157 calculated independently for each experiment.

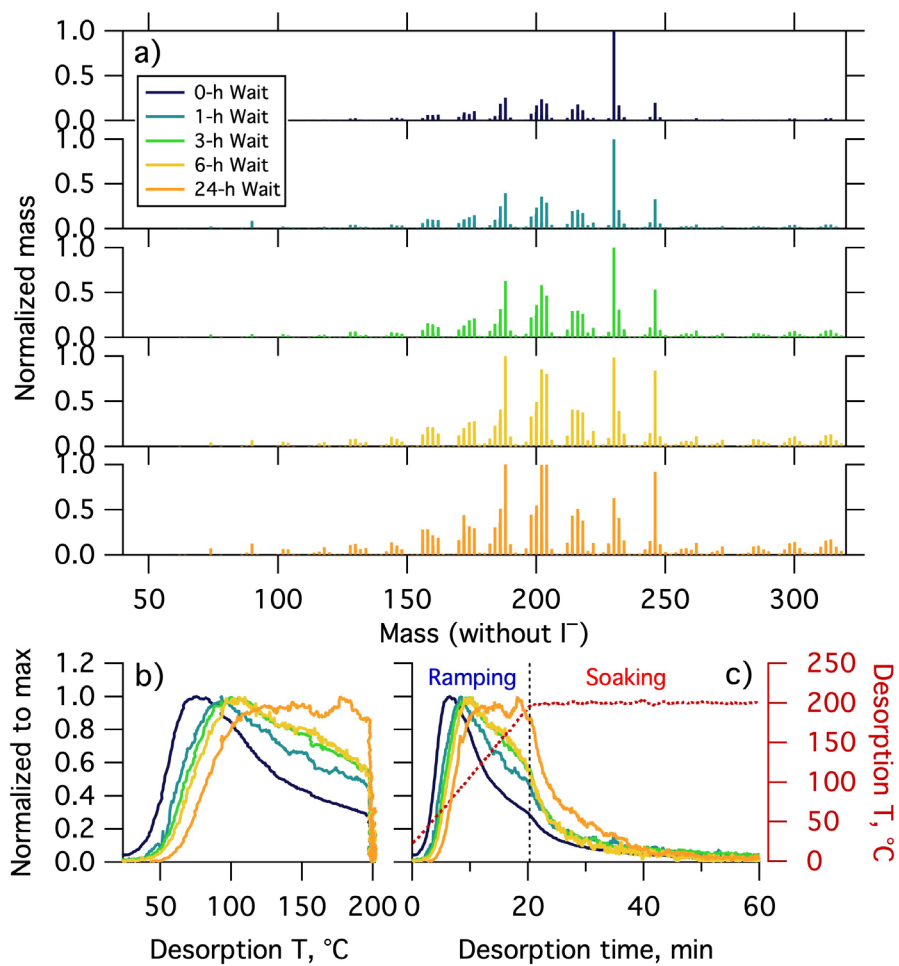
1158



1159
 1160 **Figure 10.** Multiple clustering results for α -pinene + OH SOA as a function of NO concentration. Clustering
 1161 results are separately shown for the (a) 5 ppb NO, (b) 10 ppb NO, and (c) 25 ppb NO experiments. Each
 1162 panel includes unweighted average thermograms (grey lines), mass-weighted average thermograms
 1163 (black lines) and individual cluster members (colored lines). (d) Contribution of each cluster to the total
 1164 mass for each experiment. The mass contribution of filtered-out ions (black bar) and unclustered ions
 1165 (gray bar) are also shown.

1166

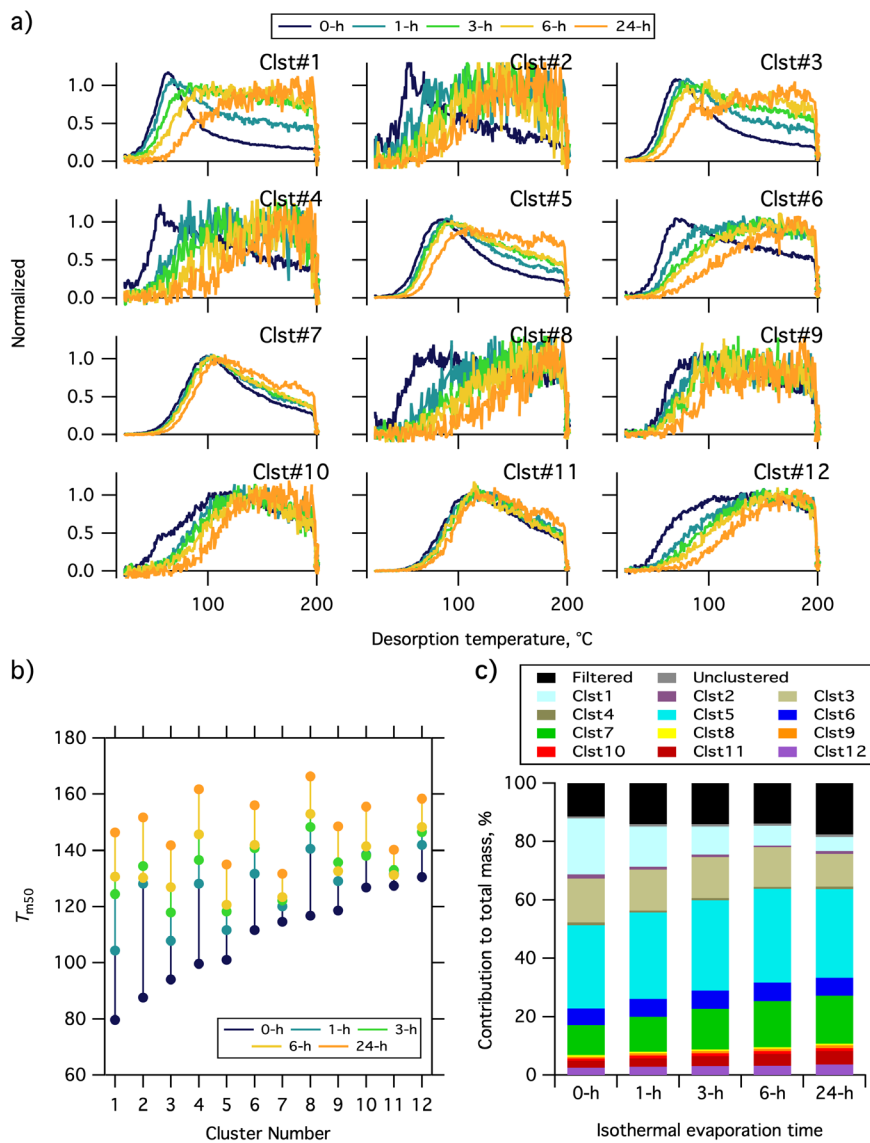
1167



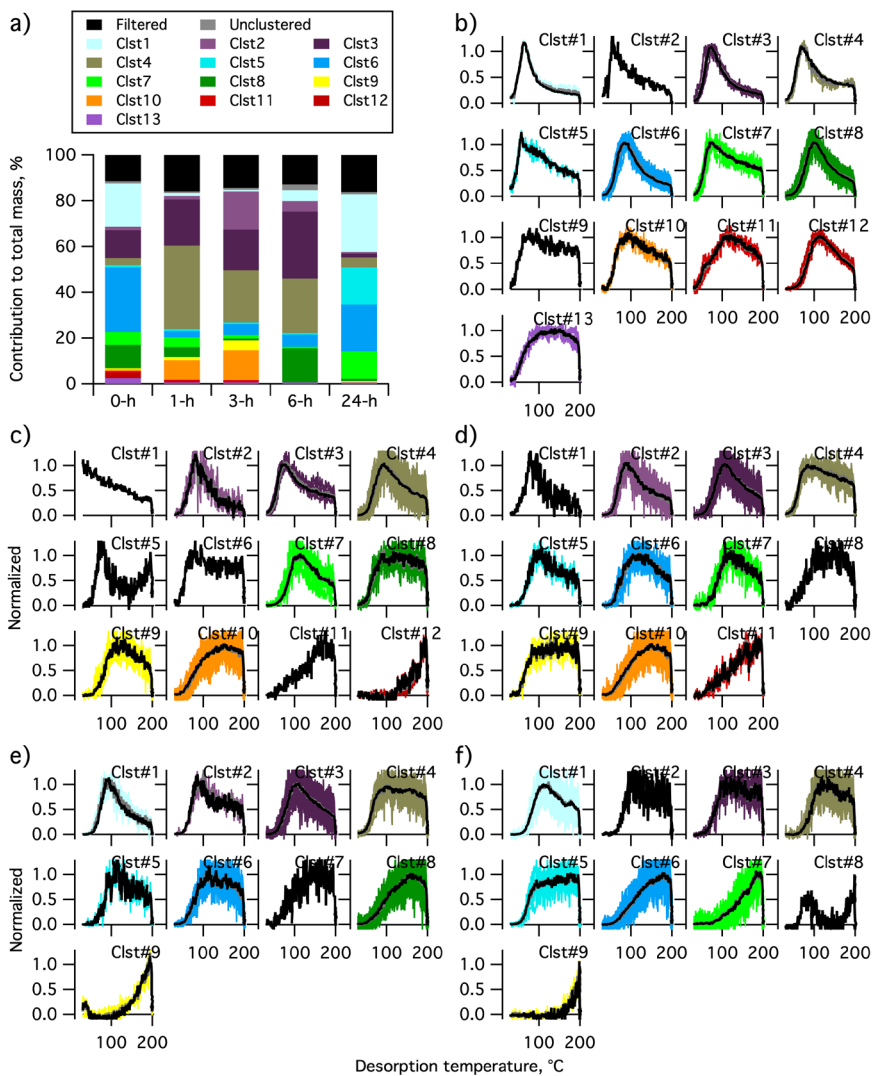
1168

1169 **Figure 11.** (a) Normalized mass spectra of α -pinene + O_3 SOA measured after different extents of
1170 isothermal evaporation at room temperature. The mass excludes iodine. The normalized thermograms of
1171 bulk SOA versus (b) temperature and (c) time, with the desorption temperature shown as a red dashed
1172 line. The vertical black dashed line in (c) delineates between ramping and soaking. The mass spectrum or
1173 thermogram colors indicate the isothermal evaporation time (see legend), with darker colors indicating
1174 shorter times.

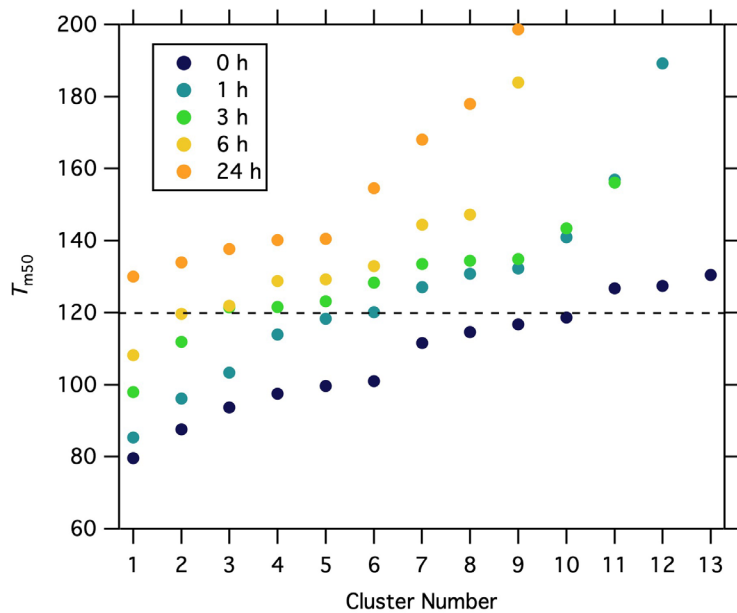
1175



1176
 1177 **Figure 12.** Single clustering results for α -pinene + O_3 SOA for different isothermal evaporation times. (a)
 1178 Comparison of the normalized, weighted-average thermograms of the 12 clusters of 0-h wait (navy), 1-h
 1179 wait (blue), 3-h wait (green), 6-h wait (yellow) and 24-h wait (orange) experiments. Note that the
 1180 absolute signals of all of the clusters decrease with evaporation, but to varying extents (**Figure S6**).



1181
 1182 **Figure 13.** Multiple clustering results for α -pinene + O₃ SOA as a function of isothermal evaporation time.
 1183 (a) Contribution of each cluster to the total mass for each experiment, along with the contributions of
 1184 filtered-out ions (black bar) and unclustered ions (gray bar). The number of clusters obtained generally
 1185 decreases with isothermal evaporation time. (b-f) The unweighted average (gray) and mass-weighted
 1186 average (black) thermograms, along with the thermograms of individual members of clusters for the (b)
 1187 0-h, (c) 1-h, (d) 3-h, (e) 6-h, and (f) 24-h wait experiments. The cluster colors are consistent between panels.



1188
 1189 **Figure 14.** The T_{m50} values of the cluster-specific thermograms from multiple clustering for the five
 1190 isothermal evaporation experiments.

1191

Supplemental Material for

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

The supplemental material includes six tables and six figures, along with additional experimental details for the experiments discussed here

FIGAERO-CIMS Instrument Description

The FIGAERO-CIMS instrument has been described previously in detail (Lee et al., 2014; Lopez-Hilfiker et al., 2014). In brief, the measurement of organic aerosol using FIGAERO-CIMS involves two steps: real-time sampling of the gas-phase with simultaneous isothermal collection of particles onto a filter through a separate inlet, followed by temperature programmed thermal desorption and detection of particle-phase species. Thermal desorption of particles occurs in two-stages: a “ramping” and “soaking” period. During ramping, the temperature of UHP N₂ programmatically increases from room temperature to 200 °C, typically at 10 °C min⁻¹. The majority of the organic aerosol mass desorbs during the ramping stage. During the soaking period, the UHP N₂ is held at 200 °C for ca. 30–40 mins to facilitate evaporation of the remaining, low-volatility organic mass from the filter.

The desorbed gas-phase compounds are transferred to the high-resolution time-of-flight (HRTof) CIMS for continuous detection and quantification at ca. 1 Hz. Iodide (I⁻) is used as the

reagent ion, which is appropriate for characterization of generally highly oxygenated components comprising most secondary organic aerosol (Isaacman-VanWertz et al., 2017; Lee et al., 2018). In a typical SOA system, hundreds of ions having the general formula $C_xH_yO_zI^-$ are usually detected. The resulting signal or mass concentration versus temperature (or equivalently time) curves for each ion constitute a thermogram. The overall bulk thermogram is obtained by summing together the individual thermograms.

All individual thermograms are background corrected by subtracting the observed thermograms from appropriate blank experiments. Blank experiments are periodically conducted by placing an additional Teflon filter upstream of the particle filter. The same collection time and desorption processes are used for blank experiments as for samples. The blanks account for contributions from adsorption of gaseous compounds in the air stream and for desorption of compounds from the inner surfaces of the FIGAERO.

Additional Experimental Details

Several example applications of the clustering on FIGAERO-CIMS data are discussed in Section 4. These include experiments on SOA derived from: (1) OH + α -pinene and (2) OH + Δ -3-carene, both at low-NO_x conditions; (3) OH + α -pinene as a function of [NO]; and (4) O₃ + α -pinene, but where the SOA is allowed to isothermally evaporate for varying amounts of time prior to thermal desorption. These experiments are briefly described below.

All the experiments were done in a 10.6 m³ Teflon environmental chamber at Pacific Northwest National Laboratory (PNNL) (Liu et al., 2012; Liu et al., 2016). Details of SOA formation and chamber conditions are summarized in **Table 1**.

Experiments #1-3 were part of the campaign of SOA Formation from Forest Emissions Experiment (SOAFFEE). SOAFFEE was designed and conducted to study the influence of reaction conditions on the formation, composition and properties of biogenic SOA. We consider only a subset of all the SOAFFEE experiments. For the experiments in this study, the chamber operated in continuous-flow mode. The total flow through the chamber was 48.2 L min⁻¹, resulting in a residence time of ~3.7 hours. Biogenic precursors were delivered into the chamber by flushing pure air through a glass bulb immersed in a temperature-controlled liquid bath held at 1 °C that

contained a small volume of the pure liquid. OH radicals were produced from the photolysis of H₂O₂. An aqueous solution of H₂O₂ was introduced into a gently warmed glass bulb by a syringe pump. A controlled flow of pure air is passed through the bulb to deliver the desired concentration of H₂O₂ into the chamber.

Seed particles of (NH₄)₂SO₄ were used to enhance SOA formation and reduce losses of semi-volatile reaction products to the chamber walls. Seed particles from atomization were dried and 50 nm particles were selected using a differential mobility analyzer (DMA), which were introduced into the chamber. The chamber relative humidity (RH) was 50%. For experiments #1-2, no NO was added. For experiment #3, a varying amount of NO was added to the chamber via a calibrated gas cylinder and a mass flow controller. A suite of online instruments characterized the chamber outflow, including a UV absorption O₃ analyzer (Thermo Environmental Instruments model 49C), an NO-NO₂-NO_x analyzer (Thermo Environmental models 42c and 42i), a TSI scanning mobility particle sizer for the number and volume concentrations of aerosols (SMPS Model 3081), an Ionicon quadrupole proton-transfer-reaction mass spectrometer (PTR-QMS) for concentration of precursors, and an Aerodyne high-resolution time-of-flight mass spectrometer (HR-ToF-AMS) for the submicron particle mass and bulk composition. Additionally, FIGAERO-CIMS was used to monitor the gas- and particle-phase products of VOC oxidation.

Experiment #4 has been described in detail previously (D'Ambro et al., 2018). The work herein focuses on the experiments performed at PNNL at an evaporation RH of 80%. For this set of experiments, FIGAERO-CIMS was operated in two modes, normal and wait mode. In the normal mode, desorption is initiated as soon as sufficient mass is collected. In the wait mode, collected particles were allowed to isothermally evaporate for some period of time prior to thermal desorption. For the isothermal evaporation, the UHP N₂ humidified to 80% RH was continuously passed over the filter at room temperature. Dilution of the air around the filter led to evaporation of SOA. The time of isothermal evaporation ranged from 1 hour to 24 hours, resulting in varying extents of mass loss of SOA from the filter. The chemical compositions of the remaining SOA were then characterized by thermal desorption of the particles.

Table S1 Chemical characteristics of each cluster identified in the α -pinene + OH SOA system

Cluster #	Expt. #1 (α -pinene + OH) Molecular Formula	O:C	H:C	MW	Mass %	# Ions	$T_{m,50}$	$T_{m,75}$	ΔT
1	C _{9.3} H _{14.9} O _{3.3} N _{0.0}	0.36	1.60	179.3	2.8	3	86.6	110.7	24.1
2	C _{9.6} H _{16.6} O _{5.1} N _{0.0}	0.53	1.72	213.4	20.1	11	93.7	120.6	26.9
3	C _{8.5} H _{13.0} O _{5.1} N _{0.0}	0.60	1.54	196.6	44.3	42	103.7	139.2	35.5
4	C _{8.2} H _{12.9} O _{6.7} N _{0.0}	0.81	1.56	218.5	8.2	13	110.6	140.8	30.2
5	C _{7.4} H _{10.8} O _{4.0} N _{0.0}	0.54	1.46	163.6	4.0	3	116.2	153.1	36.9
6	C _{11.1} H _{17.8} O _{8.1} N _{0.0}	0.73	1.60	280.6	0.3	6	126.2	154.5	28.3
7	C _{15.3} H _{22.6} O _{5.4} N _{0.0}	0.35	1.48	292.6	0.2	3	129.5	162.4	32.9
8	C _{8.5} H _{10.9} O _{4.0} N _{0.0}	0.47	1.29	176.9	1.8	2	131.0	164.1	33.1
9	C _{7.2} H _{9.4} O _{4.5} N _{0.0}	0.62	1.31	167.8	9.8	17	131.3	163.7	32.4
10	C _{6.9} H _{7.9} O _{3.5} N _{0.0}	0.50	1.15	146.7	0.6	4	145.7	172.1	26.4
11	C _{4.1} H _{4.0} O _{3.9} N _{0.0}	0.93	0.97	115.6	0.3	2	161.8	184.2	22.4
Unclustered					0	0			
Filtered					7.5	188			

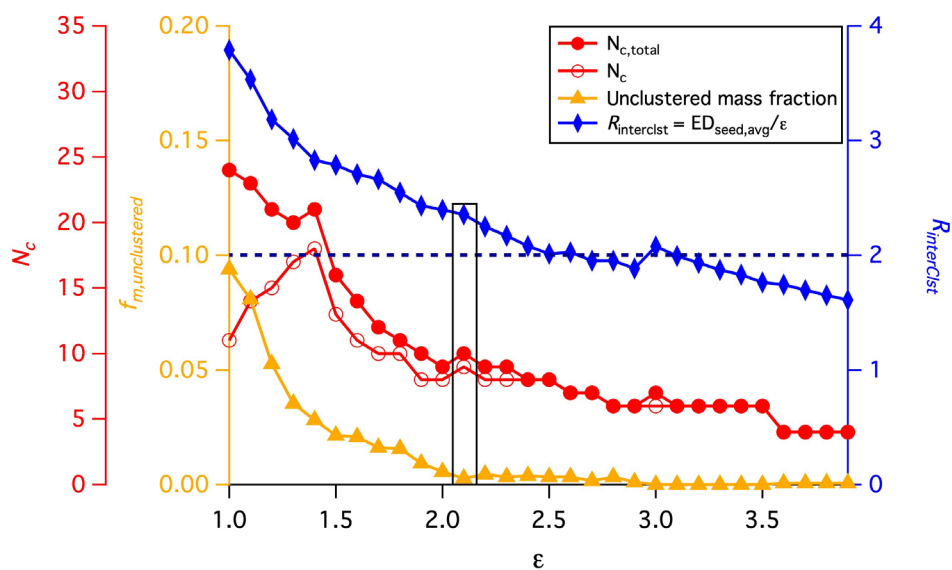


Figure S1. Similar to Figure 4, but guidance for determining the optimal ϵ for the Δ -3-carene +OH SOA. The variation of four parameters, N_c , $N_{c,total}$, $f_{m,unclustered}$ and $R_{interClst}$ are shown as a function of the distance criterion ϵ . The black horizontal dashed line is guide judgement for $R_{interClst} \geq 2$. The values highlighted by a rectangle are those corresponding to the optimal ϵ used for the following clustering analysis.

Table S2 Chemical characteristics of each cluster identified in the Δ -3-carene + OH SOA system

Cluster #	Expt. #2 (Δ -3-carene+ OH) Molecular Formula	O:C	H:C	MW	Mass %	# Ions	$T_{m,50}$	$T_{m,75}$	ΔT
1	C _{9.5} H _{16.7} O _{5.7} N _{0.0}	0.60	1.76	221.9	16.6	3	99.4	127.0	27.6
2	C _{6.9} H _{11.3} O _{4.7} N _{0.1}	0.68	1.64	170.7	3.4	11	102.3	138.7	36.4
3	C _{8.3} H _{12.9} O _{5.4} N _{0.0}	0.65	1.55	198.9	17.5	42	106.9	128.3	21.4
4	C _{8.0} H _{12.4} O _{5.3} N _{0.0}	0.65	1.54	193.2	28.0	13	110.1	143.9	33.8
5	C _{11.4} H _{19.7} O _{8.6} N _{0.0}	0.76	1.73	294.1	0.4	3	120.1	147.0	26.9
6	C _{8.2} H _{11.7} O _{5.2} N _{0.0}	0.63	1.42	193.3	14.8	6	121.4	156.1	34.7
7	C _{6.6} H _{8.6} O _{4.6} N _{0.0}	0.69	1.30	161.4	6.9	3	131.3	164.6	33.3
8	C _{6.9} H _{7.9} O _{3.9} N _{0.0}	0.56	1.13	153.1	1.8	2	141.4	170.1	28.7
9	C _{6.0} H _{7.2} O _{4.0} N _{0.0}	0.68	1.21	143.2	0.3	17	151.2	176.3	25.1
10	C _{4.0} H _{4.0} O _{4.0} N _{0.0}	1.00	1.00	116	0.6	4	157.5	181.9	24.4
Unclustered					0.3	4			
Filtered					9.3	183			

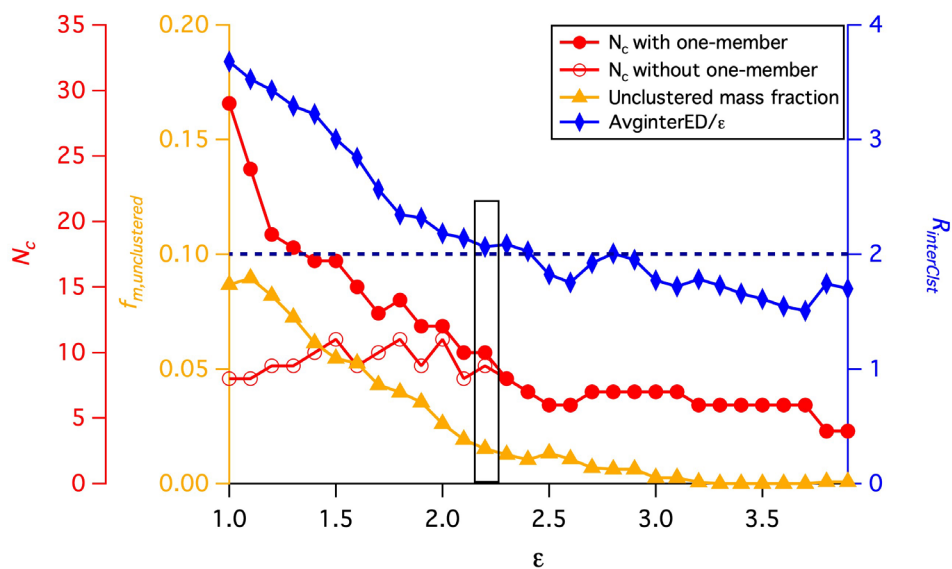


Figure S2. Similar to Figure 4, but guidance for determining the optimal ϵ for the α -pinene +OH SOA formed under 5 ppb NO for the single clustering approach. The variation of four parameters, N_c , $N_{c,total}$, $f_{m,unclustered}$ and $R_{interCist}$ are shown as a function of the distance criterion ϵ . The black horizontal dashed line is guide judgement for $R_{interCist} \geq 2$. The values highlighted by a rectangle are those corresponding to the optimal ϵ used for the following clustering analysis.

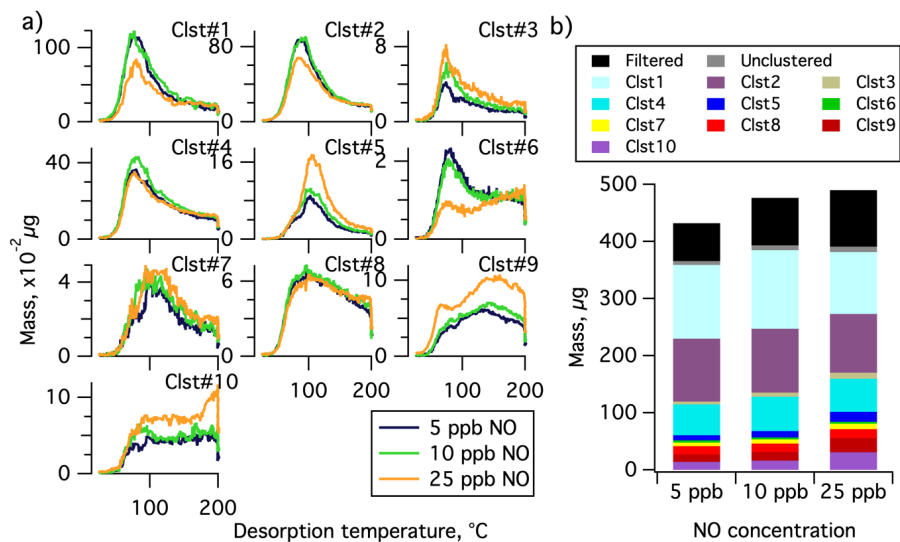


Figure S3 Similar to **Figure 9** but presented in the absolute sense. (a) Comparison of the summed thermograms of the 10 clusters of 5 ppb (navy), 10 ppb (green) and 25 ppb (orange) NO experiments. (b) Absolute mass of each cluster for each experiment, including the summed mass of filtered out ions (black) and unclustered ions (gray).

Table S3. Chemical characteristics of each cluster identified in the α -pinene + OH + NO SOA system. The single clustering approach is used based on 5 ppb NO experiment.

Table S3-1.

Cluster #	Expt. #3 (α -pinene + OH + NO)			[NO]	[NO]	[NO]	# Ions	
	Molecular Formula	O:C	H:C	5 ppb	10 ppb	25 ppb		
1	C _{8.8} H _{14.3} O _{4.7} N _{0.0}	0.53	1.62	195.1	29.9	28.8	27	
2	C _{9.3} H _{14.8} O _{6.0} N _{0.0}	0.65	1.59	222.4	25.4	23.6	19	
3	C _{7.5} H _{10.7} O _{6.4} N _{0.4}	0.85	1.42	208.7	1.2	1.4	2	
4	C _{8.0} H _{11.7} O _{4.9} N _{0.1}	0.61	1.47	187.5	12.4	12.6	10	
5	C _{8.0} H _{12.0} O _{6.6} N _{0.0}	0.75	1.50	204.0	2.2	2.4	1	
6	C _{5.5} H _{7.1} O _{3.7} N _{0.0}	0.67	1.28	132.3	0.9	0.7	2	
7	C _{8.9} H _{11.4} O _{6.3} N _{0.0}	0.72	1.28	219.0	1.4	1.5	4	
8	C _{8.1} H _{11.1} O _{4.1} N _{0.0}	0.50	1.37	173.9	3.3	3.1	7	
9	C _{4.6} H _{5.6} O _{3.9} N _{0.0}	0.83	1.20	123.2	3.0	3.1	6	
10	C _{8.2} H _{10.2} O _{5.0} N _{0.0}	0.61	1.24	188.6	3.3	3.4	2	
Unclustered					1.5	1.7	2.0	8
Filtered					15.5	17.5	20.2	206

Table S3-2.

Cluster #	Expt. #3 (α -pinene + OH + NO) Molecular Formula	[NO]	[NO]	[NO]	[NO]	[NO]	[NO]	[NO]	[NO]	[NO]
		5 ppb	10 ppb $T_{m,50}$	25 ppb	5 ppb	10 ppb $T_{m,75}$	25 ppb	5 ppb	10 ppb ΔT	25 ppb
1	C _{8.8} H _{14.3} O _{4.7} N _{0.0}	92.4	95.7	101.9	130.0	133.1	151.1	37.6	37.4	49.2
2	C _{9.3} H _{14.8} O _{6.0} N _{0.0}	100.0	100.9	103.9	136.5	136.5	145.8	36.5	35.6	41.9
3	C _{7.5} H _{10.7} O _{6.4} N _{0.4}	101.8	97.1	101.8	145.9	142.2	146.1	44.1	45.1	44.3
4	C _{8.0} H _{11.7} O _{4.9} N _{0.1}	102.6	102.0	107.0	145.1	144.5	153.6	42.5	42.5	46.6
5	C _{8.0} H _{12.0} O _{6.6} N _{0.0}	108.4	109.4	114.0	138.0	137.3	139.3	29.6	27.9	25.3
6	C _{5.5} H _{7.1} O _{3.7} N _{0.0}	111.1	119.9	144.3	158.2	161.3	179.0	47.1	41.4	34.7
7	C _{8.9} H _{11.4} O _{6.3} N _{0.0}	122.0	119.7	125.6	156.1	152.8	159.4	34.1	33.1	33.8
8	C _{8.1} H _{11.1} O _{4.1} N _{0.0}	125.7	125.4	129.9	160.3	160.8	166.3	34.6	35.4	36.4
9	C _{4.6} H _{5.6} O _{3.9} N _{0.0}	135.9	137.4	138.9	166.4	167.1	169.3	30.5	29.7	30.4
10	C _{8.2} H _{10.2} O _{5.0} N _{0.0}	140.6	137.2	146.5	175.7	173.7	184.5	35.1	36.5	38.0
Unclustered										
Filtered										

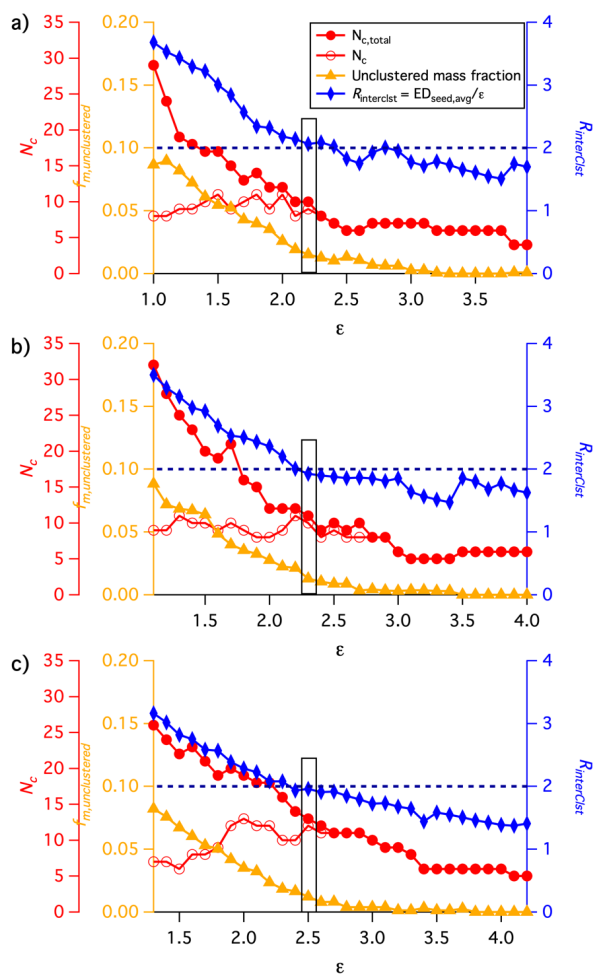


Figure S4. Guidance for determining the optimal ϵ for the α -pinene +OH SOA formed under (a) 5 ppb, (b) 10 ppb and (c) 25 ppb NO conditions for the multiple clustering approach. The variation of four parameters, N_c , $N_{c,all}$, $f_{m,unclustered}$ and $R_{interClst}$ are shown as a function of the distance criterion ϵ . The black horizontal dashed line is guide judgement for $R_{interClst} \geq 2$. The values highlighted by a rectangle are the values corresponding to the optimal ϵ used for the following clustering analysis.

Table S4 Chemical characteristics of each cluster identified in the α -pinene + OH + NO SOA system for three different NO conditions (5, 10 and 25 ppb) from the multiple-clustering approach.

Table S4-1. Results for the 5 ppb NO experiment.

Cluster #	Expt. #3a (α -pinene+ OH+5 ppb NO)							
	Molecular Formula	O:C	H:C	Mass %	# Ions	$T_{m,50}$	$T_{m,75}$	ΔT
1	C _{8.8} H _{14.3} O _{4.7} N _{0.0}	0.53	1.62	29.9	27	92.4	130.0	37.6
2	C _{9.3} H _{14.8} O _{6.0} N _{0.0}	0.65	1.59	25.4	19	100.0	136.5	36.5
3	C _{7.5} H _{10.7} O _{6.4} N _{0.4}	0.85	1.42	1.2	2	101.8	145.9	44.1
4	C _{8.0} H _{11.7} O _{4.9} N _{0.1}	0.61	1.47	12.4	10	102.6	145.1	42.5
5	C _{8.0} H _{12.0} O _{6.0} N _{0.0}	0.75	1.50	2.2	1	108.4	138.0	29.6
6	C _{5.5} H _{7.1} O _{3.7} N _{0.0}	0.68	1.28	0.9	2	111.1	158.2	47.1
7	C _{8.9} H _{11.4} O _{6.3} N _{0.0}	0.72	1.28	1.4	4	122.0	156.1	34.1
8	C _{8.1} H _{11.1} O _{4.1} N _{0.0}	0.50	1.37	3.3	7	125.7	160.3	34.6
9	C _{4.6} H _{5.6} O _{3.9} N _{0.0}	0.83	1.20	3.0	6	135.9	166.4	30.5
10	C _{8.2} H _{10.2} O _{5.0} N _{0.0}	0.61	1.24	3.3	2	140.6	175.7	35.1
Unclustered				1.5	8			
Filtered				15.5	208			

Table S4-2. Results for the 10 ppb NO experiment.

Cluster #	Expt. #3b (α -pinene+ OH+10 ppb NO)							
	Molecular Formula	O:C	H:C	Mass %	# Ions	$T_{m,50}$	$T_{m,75}$	ΔT
1	C _{8.5} H _{13.6} O _{4.9} N _{0.1}	0.58	1.60	27.6	28	95.1	134.1	39.0
2	C _{9.1} H _{14.3} O _{6.1} N _{0.1}	0.66	1.57	25.7	21	101.9	136.5	34.6
3	C _{8.5} H _{12.8} O _{4.9} N _{0.1}	0.58	1.50	17.3	14	102.1	143.3	41.2
4	C _{5.4} H _{11.1} O _{7.6} N _{0.2}	1.40	2.06	2.1	4	104.1	148.4	44.3
5	C _{6.8} H _{8.4} O _{4.6} N _{0.0}	0.68	1.23	3.2	5	110.9	159.5	48.6
6	C _{12.5} H _{18.9} O _{7.5} N _{0.0}	0.60	1.52	0.2	2	124.5	152.9	28.4
7	C _{9.7} H _{14.3} O _{7.3} N _{0.0}	0.75	1.47	0.3	2	126.5	159.3	32.8
8	C _{8.2} H _{11.0} O _{4.2} N _{0.0}	0.51	1.35	3.1	7	126.3	161.1	34.8
9	C _{4.0} H _{4.0} O _{6.0} N _{0.0}	1.50	1.00	0.5	1	134.3	168.5	34.2
10	C _{7.8} H _{9.5} O _{4.8} N _{0.0}	0.61	1.23	4.1	2	137.8	172.4	34.6
11	C _{3.2} H _{3.8} O _{3.9} N _{0.0}	1.25	1.22	1.9	4	138.6	166.3	27.7
Unclustered				1.3	10			
Filtered				12.6	195			

Table S4-3. Results for the 25 ppb NO experiment.

Cluster #	Expt. #3c (α -pinene+ OH+25 ppb NO) Molecular Formula	O:C	H:C	Mass %	# Ions	$T_{m,50}$	$T_{m,75}$	ΔT
1	C _{10.0} H _{15.0} O _{6.0} N _{1.0}	0.60	1.50	1.1	1	92.3	139.8	47.5
2	C _{6.0} H _{8.8} O _{5.0} N _{0.0}	0.82	1.46	0.8	4	98.5	147.8	49.3
3	C _{8.1} H _{12.4} O _{5.4} N _{0.2}	0.66	1.53	39.8	29	103.4	150.7	47.3
4	C _{9.1} H _{14.1} O _{6.3} N _{0.1}	0.69	1.55	13.7	18	106.5	144.4	37.9
5	C _{6.1} H _{9.8} O _{6.3} N _{0.0}	1.03	1.60	3.7	2	121.2	160.9	39.7
6	C _{8.7} H _{11.3} O _{6.2} N _{0.0}	0.72	1.30	3.2	5	122.2	155.0	32.8
7	C _{7.9} H _{10.7} O _{4.8} N _{0.1}	0.61	1.36	2.8	8	125.1	165.1	40.0
8	C _{4.0} H _{5.6} O _{3.9} N _{0.0}	0.97	1.37	1.7	5	134.9	168.6	33.7
9	C _{11.8} H _{18.5} O _{7.7} N _{0.0}	0.65	1.57	0.2	3	136.5	166.2	29.7
10	C _{2.5} H _{3.0} O _{4.0} N _{0.0}	1.58	1.20	2.3	3	139.2	168.7	29.5
11	C _{5.2} H _{5.8} O _{5.1} N _{0.0}	0.97	1.11	3.8	3	140.2	175.5	35.3
12	C _{7.8} H _{9.7} O _{4.9} N _{0.0}	0.63	1.25	1.0	2	147.7	185.4	37.7
13	C _{4.2} H _{4.0} O _{3.8} N _{0.0}	0.92	0.96	1.5	2	167.5	188.9	21.4
Unclustered				0.7	7			
Filtered				16.0	205			

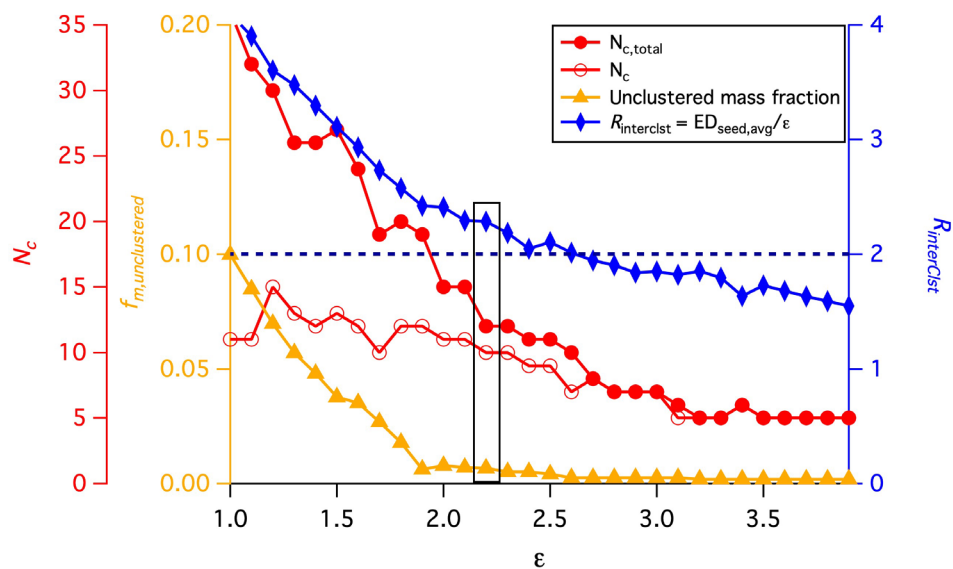


Figure S5. Guidance for determining the optimal ϵ for the α -pinene + O_3 SOA system with no isothermal evaporation for the single clustering approach. The variation of four parameters, N_c , $N_{c,total}$, $f_{m,unclustered}$ and $R_{interClst}$ are shown as a function of the distance criterion ϵ . The black horizontal dashed line is guide judgement for $R_{interClst} \geq 2$. The values highlighted by a rectangle are the values corresponding to the optimal ϵ used for the following clustering analysis.

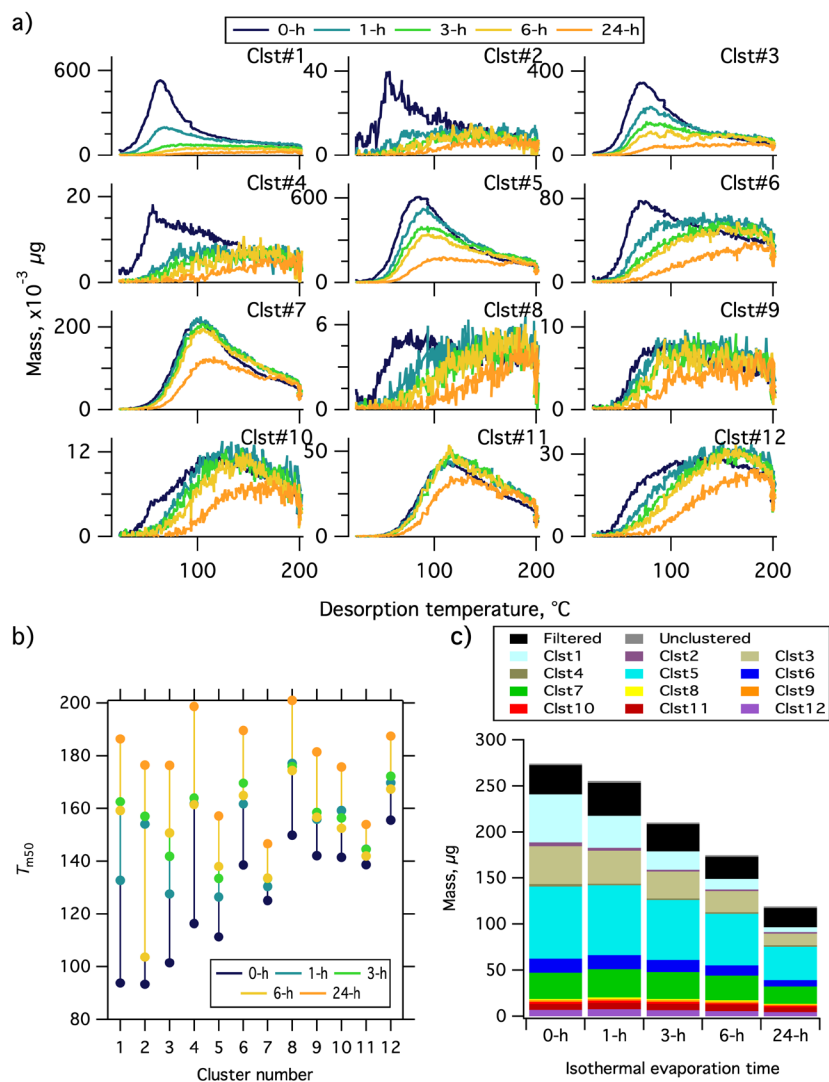


Figure S6 Similar to **Figure 12** but presented in the absolute sense. (a) Comparison of the summed thermograms of the 12 clusters of 0-h wait (navy), 1-h wait (blue), 3-h wait (green), 6-h wait (yellow) and 24-h wait (orange) experiments. (b) Changes in the T_{m50} for all the clusters calculated from the summed thermograms, with the same color scheme as (a). (c) Absolute mass of each cluster for each experiment, including the summed mass of filtered out ions (black) and unclustered ions (gray).

Table S5 Chemical characteristics of each cluster identified in the α -pinene + O₃ + Evaporation SOA system from the single-clustering approach, using the no-wait experiment as the reference.

Table S5-1.

Cluster #	Expt. #4 (α -pinene+ O ₃ +Evap.) Molecular Formula	0 h 1 h 3 h 6 h 24 h								
		O:C	H:C	MW	Mass %					# Ions
1	C _{9.8} H _{14.0} O _{5.7} N _{0.0}	0.58	1.42	222.8	19.0	13.7	9.5	7.0	4.8	2
2	C _{9.0} H _{16.0} O _{3.0} N _{0.0}	0.33	1.78	172.0	1.5	1.1	0.9	0.3	1.0	1
3	C _{8.7} H _{14.1} O _{4.8} N _{0.0}	0.54	1.61	195.3	15.0	14.0	14.0	13.5	11.3	9
4	C _{6.5} H _{11.1} O _{4.0} N _{0.0}	0.61	1.69	153.1	1.0	0.7	0.7	0.7	0.8	2
5	C _{8.9} H _{14.0} O _{6.1} N _{0.0}	0.69	1.58	218.4	28.5	29.6	31.0	32.1	30.3	35
6	C _{8.7} H _{12.2} O _{4.9} N _{0.0}	0.56	1.40	195.0	5.6	6.1	6.2	6.3	6.2	8
7	C _{12.6} H _{23.5} O _{7.8} N _{0.0}	0.62	1.86	299.5	10.4	12.1	14.1	15.8	16.4	34
8	C _{6.0} H _{8.0} O _{4.0} N _{0.0}	0.67	1.33	144.0	0.46	0.5	0.4	0.5	0.6	1
9	C _{6.7} H _{13.9} O _{8.2} N _{0.0}	1.22	2.07	225.5	0.6	0.7	0.7	0.8	0.9	3
10	C _{5.1} H _{8.1} O _{3.0} N _{0.0}	0.59	1.60	117.3	0.8	1.0	1.0	0.9	1.0	4
11	C _{13.7} H _{24.7} O _{8.5} N _{0.0}	0.62	1.80	325.1	2.4	2.8	3.4	4.0	4.6	8
12	C _{8.1} H _{10.8} O _{4.1} N _{0.0}	0.50	1.34	173.6	2.5	2.9	3.1	3.2	3.6	5
Unclustered					0.7	0.8	0.8	0.8	0.8	5
Filtered					11.5	14.0	14.0	13.8	17.6	185

Table S5-2

Cluster #	Expt. #4 (α -pinene+ O ₃ +Evap.) Molecular Formula	0 h 1 h 3 h 6 h 24 h 0 h 1 h 3 h 6 h 24 h									
		<i>T</i> _{m,50}					ΔT				
1	C _{9.8} H _{14.0} O _{5.7} N _{0.0}	79.6	104.3	124.4	130.6	146.4	40.1	45.9	38.3	33.4	30.1
2	C _{9.0} H _{16.0} O _{3.0} N _{0.0}	87.5	128.1	134.3	130.2	151.7	43.2	35.6	30.3	23.8	24.6
3	C _{8.7} H _{14.1} O _{4.8} N _{0.0}	93.9	107.7	117.8	126.9	141.7	33.7	40.6	38.6	34.6	32.0
4	C _{6.5} H _{11.1} O _{4.0} N _{0.0}	99.6	128.1	136.5	145.6	161.7	40.9	36.7	31.6	25.4	21.9
5	C _{8.9} H _{14.0} O _{6.1} N _{0.0}	101.0	111.5	118.3	120.6	135.0	32.5	34.5	35.8	35.4	32.6
6	C _{8.7} H _{12.2} O _{4.9} N _{0.0}	111.6	131.7	140.8	141.8	156.0	42.2	33.4	29.2	28.7	25.3
7	C _{12.6} H _{23.5} O _{7.8} N _{0.0}	114.6	120.0	122.1	123.3	131.6	32.0	32.0	31.7	32.4	30.5
8	C _{6.0} H _{8.0} O _{4.0} N _{0.0}	116.7	140.5	148.3	153.0	166.4	40.9	32.6	27.4	23.3	19.5
9	C _{6.7} H _{13.9} O _{8.2} N _{0.0}	118.6	129.0	135.7	132.6	148.5	37.6	34.3	28.8	32.0	30.4
10	C _{5.1} H _{8.1} O _{3.0} N _{0.0}	126.8	138.5	138.0	141.4	155.5	31.4	28.0	28.3	23.3	21.1
11	C _{13.7} H _{24.7} O _{8.5} N _{0.0}	127.4	131.6	133.0	131.2	140.2	29.0	29.0	28.0	27.8	26.4
12	C _{8.1} H _{10.8} O _{4.1} N _{0.0}	130.5	141.9	146.4	148.3	158.4	33.0	29.7	27.4	25.1	23.2
Unclustered											
Filtered											

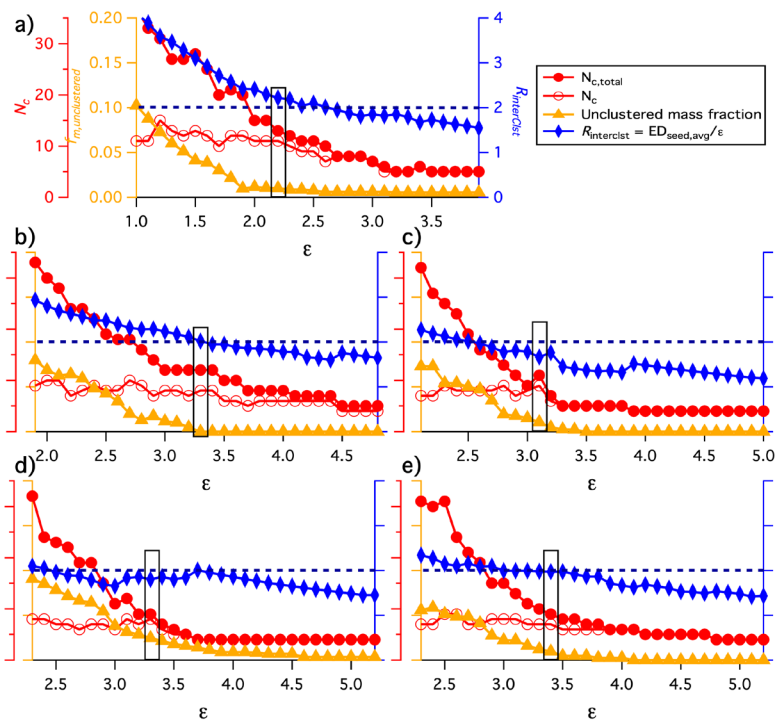


Figure S7. Guidance for determining the optimal ϵ for the α -pinene + O_3 SOA system with (a) no isothermal evaporation (b) 1 hr (c) 3 hrs (d) 6 hrs and (e) 24 hrs of isothermal evaporation for the multiple clustering approach. The variation of four parameters, N_c , $N_{c,total}$, $f_{m,unclustered}$ and $R_{interClist}$ are shown as a function of the distance criterion ϵ . The black horizontal dashed line is guide judgement for $R_{interClist} \geq 2$. The values highlighted by a rectangle are the values corresponding to the optimal ϵ used for the following clustering analysis.

Table S6 Chemical characteristics of each cluster identified in the α -pinene + O₃ + Evaporation SOA system for five different isothermal evaporation conditions. The multiple clustering approach is used so each evaporation experiment is clustered independently.

Table S6-1. Results for the 0 h isothermal evaporation experiment.

Cluster #	Expt. #4a (α -pinene+ O ₃ + 0-h Evap.)							
	Molecular Formula	O:C	H:C	Mass %	# Ions	T _{m,50}	T _{m,75}	ΔT
1	C _{9,8} H _{14,0} O _{5,7} N _{0,0}	0.58	1.42	18.9	2	79.6	119.7	40.1
2	C _{9,0} H _{16,0} O _{3,0} N _{0,0}	0.33	1.78	1.4	1	87.5	130.8	43.3
3	C _{8,5} H _{13,6} O _{4,9} N _{0,0}	0.58	1.61	12.2	8	93.6	124.3	30.7
4	C _{10,9} H _{18,5} O _{3,7} N _{0,0}	0.34	1.69	3.2	2	97.5	140.7	43.2
5	C _{6,5} H _{11,1} O _{4,0} N _{0,0}	0.61	1.69	1.0	2	99.6	140.5	40.9
6	C _{8,9} H _{14,0} O _{6,1} N _{0,0}	0.69	1.58	28.2	35	101.0	133.5	32.5
7	C _{8,7} H _{12,2} O _{4,9} N _{0,0}	0.56	1.40	5.6	8	111.6	153.8	42.2
8	C _{12,6} H _{23,5} O _{7,8} N _{0,0}	0.62	1.86	10.3	34	114.6	146.5	31.9
9	C _{6,0} H _{8,0} O _{4,0} N _{0,0}	0.67	1.33	0.5	1	116.7	157.6	40.9
10	C _{6,7} H _{13,9} O _{8,2} N _{0,0}	1.22	2.07	0.6	3	118.6	156.2	37.6
11	C _{5,1} H _{8,1} O _{3,0} N _{0,0}	0.59	1.60	0.8	4	126.8	158.2	31.4
12	C _{13,7} H _{24,7} O _{8,5} N _{0,0}	0.62	1.80	2.3	8	127.4	156.3	28.9
13	C _{8,1} H _{10,8} O _{4,1} N _{0,0}	0.50	1.34	2.5	5	130.5	163.5	33.0
Unclustered				1.0	7			
Filtered				11.4	191			

Table S6-2. Results for the 1 h isothermal evaporation experiment.

Cluster #	Expt. #4b (α -pinene+ O ₃ + 1-h Evap.)							
	Molecular Formula	O:C	H:C	Mass %	# Ions	T _{m,50}	T _{m,75}	ΔT
1	C _{2,0} H _{4,0} O _{3,0} N _{0,0}	1.50	2.00	1.3	1	85.4	131.6	46.2
2	C _{9,4} H _{16,7} O _{5,2} N _{0,0}	0.55	1.79	1.6	3	96.1	125.8	29.7
3	C _{9,1} H _{13,5} O _{5,5} N _{0,0}	0.61	1.48	20.1	6	103.3	146.5	43.2
4	C _{9,6} H _{15,7} O _{6,2} N _{0,0}	0.65	1.63	36.5	48	113.9	148.4	34.5
5	C _{16,0} H _{32,0} O _{2,0} N _{0,0}	0.12	2.00	0.9	1	118.3	181.6	63.3
6	C _{9,0} H _{14,0} O _{4,0} N _{0,0}	0.444	1.56	2.6	1	120.0	162.0	42.0
7	C _{14,5} H _{28,1} O _{8,7} N _{0,0}	0.60	1.94	4.3	10	127.1	155.5	28.4
8	C _{8,5} H _{12,5} O _{5,4} N _{0,0}	0.64	1.47	4.3	6	130.8	164.3	33.5
9	C _{11,1} H _{20,9} O _{8,9} N _{0,0}	0.80	1.88	1.4	5	132.3	161.8	29.5
10	C _{7,4} H _{10,1} O _{3,9} N _{0,0}	0.52	1.36	8.6	18	140.9	171.0	30.1
11	C _{7,0} H _{8,0} O _{4,0} N _{0,0}	0.57	1.14	0.6	1	156.8	181.4	24.6
12	C _{17,3} H _{34,0} O _{2,1} N _{0,0}	0.12	1.97	1.1	2	189.2	199.1	9.0
Unclustered				0.0	0			
Filtered				16.5	210			

Table S6-3. Results for the 3 h isothermal evaporation experiment.

Cluster #	Expt. #4c (α -pinene+ O ₃ + 3-h Evap.)							
	Molecular Formula	O:C	H:C	Mass %	# Ions	T _{m,50}	T _{m,75}	ΔT
1	C _{10.0} H _{18.0} O _{5.0} N _{0.0}	0.50	1.80	0.5	1	97.9	133.4	35.5
2	C _{9.4} H _{14.9} O _{6.0} N _{0.0}	0.63	1.58	16.7	13	111.8	146.2	34.4
3	C _{13.1} H _{23.8} O _{7.5} N _{0.0}	0.57	1.81	17.7	31	121.4	153.3	31.9
4	C _{8.6} H _{12.8} O _{5.3} N _{0.0}	0.62	1.49	23.1	9	121.6	159.0	37.4
5	C _{15.5} H _{26.8} O _{5.9} N _{0.0}	0.38	1.73	0.4	2	123.1	156.6	33.5
6	C _{8.4} H _{12.9} O _{6.9} N _{0.0}	0.81	1.53	4.8	8	128.3	159.0	30.7
7	C _{11.1} H _{21.9} O _{8.8} N _{0.0}	0.79	1.97	1.5	5	133.4	162.1	28.7
8	C _{9.0} H _{16.0} O _{3.0} N _{0.0}	0.33	1.78	0.9	1	134.3	164.7	30.4
9	C _{7.1} H _{12.8} O _{5.7} N _{0.0}	0.81	1.80	4.3	3	134.9	169.3	34.4
10	C _{7.9} H _{10.8} O _{4.2} N _{0.0}	0.53	1.36	13.0	22	143.4	171.6	28.2
11	C _{8.3} H _{10.6} O _{4.0} N _{0.0}	0.48	1.28	1.7	2	156.1	181.3	25.2
Unclustered				1.1	5			
Filtered				14.3	205			

Table S6-4. Results for the 6 h isothermal evaporation experiment.

Cluster #	Expt. #4d (α -pinene+ O ₃ + 6-h Evap.)							
	Molecular Formula	O:C	H:C	Mass %	# Ions	T _{m,50}	T _{m,75}	ΔT
1	C _{9.8} H _{14.1} O _{6.9} N _{0.0}	0.70	1.44	4.9	3	108.2	138.3	30.1
2	C _{9.3} H _{14.8} O _{5.1} N _{0.0}	0.55	1.59	4.5	2	119.7	157.1	37.4
3	C _{11.9} H _{21.0} O _{7.3} N _{0.0}	0.61	1.76	29.2	36	121.9	154.2	32.3
4	C _{8.5} H _{13.1} O _{5.4} N _{0.0}	0.63	1.54	24.0	13	128.8	163.3	34.5
5	C _{11.5} H _{23.0} O _{9.5} N _{0.0}	0.83	2.00	0.6	2	129.2	161.7	32.5
6	C _{8.1} H _{12.7} O _{6.6} N _{0.0}	0.82	1.58	5.5	11	132.9	164.9	32.0
7	C _{11.0} H _{16.0} O _{5.0} N _{0.0}	0.45	1.45	0.4	1	144.4	172.3	27.9
8	C _{7.7} H _{10.6} O _{4.0} N _{0.0}	0.52	1.39	15.0	25	147.2	173.0	25.8
9	C _{2.7} H _{5.3} O _{3.3} N _{0.0}	1.25	2.00	0.6	2	183.9	194.9	11.0
Unclustered				2.5	11			
Filtered				12.8	203			

Table S6-5. Results for the 24 h isothermal evaporation experiment.

Cluster #	Expt. #4e (α -pinene+O ₃ +24-h Evap.)							
	Molecular Formula	O:C	H:C	Mass %	# Ions	T _{m,50}	T _{m,75}	ΔT
1	C _{12.1} H _{21.2} O _{7.2} N _{0.0}	0.60	1.75	25.4	26	130.0	161.1	31.1
2	C _{7.0} H _{16.0} O _{9.0} N _{0.0}	1.29	2.29	0.6	1	133.9	166.7	32.8
3	C _{11.3} H _{17.9} O _{6.7} N _{0.0}	0.59	1.58	1.7	4	137.7	167.7	30.0
4	C _{11.9} H _{20.0} O _{8.2} N _{0.0}	0.69	1.68	4.4	10	140.1	167.6	27.5
5	C _{8.9} H _{13.8} O _{5.5} N _{0.0}	0.62	1.54	16.1	8	140.4	171.5	31.1
6	C _{7.9} H _{12.0} O _{4.6} N _{0.0}	0.58	1.51	20.3	23	154.5	179.0	24.5
7	C _{5.2} H _{7.1} O _{3.4} N _{0.0}	0.66	1.37	12.2	17	168.0	186.3	18.3
8	C _{16.0} H _{32.0} O _{2.0} N _{0.0}	0.12	2.00	1.0	1	177.9	197.9	20.0
9	C _{17.5} H _{35.0} O _{2.1} N _{0.0}	0.12	2.00	1.2	2	198.7	199.3	0.6
Unclustered				1.0	4			
Filtered				16.1	213			

Supplemental References

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

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

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

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

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

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

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