

1 A new processing scheme for ultra-high resolution direct infusion 2 mass spectrometry data

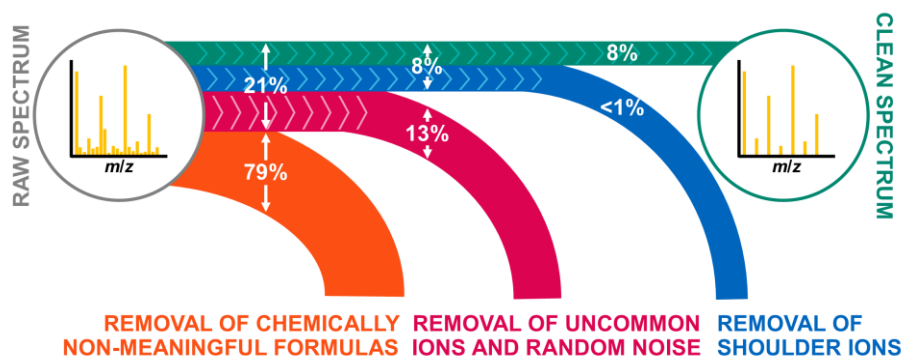
3 Arthur T. Zielinski¹, Ivan Kourtchev¹, Claudio Bortolini², Stephen J. Fuller¹, Chiara Giorio¹, Olalekan
4 A. M. Popoola¹, Sara Bogialli², Andrea Tapparò², Roderic L. Jones¹, and Markus Kalberer^{1*}

5
6 ¹ Department of Chemistry, University of Cambridge, Lensfield Road, Cambridge CB2 1EW, United Kingdom

7 ² Dipartimento di Scienze Chimiche, Università degli Studi di Padova, via Marzolo 1, 35131 Padova, Italy

8 Abstract

9 High resolution, high accuracy mass spectrometry is widely used to characterise environmental
10 or biological samples with highly complex composition enabling the identification of chemical
11 composition of often unknown compounds. Despite instrumental advancements, the accurate
12 molecular assignment of compounds acquired in high resolution mass spectra remains time
13 consuming and requires automated algorithms, especially for samples covering a wide mass
14 range and large numbers of compounds. A new processing scheme is introduced implementing
15 filtering methods based on element assignment, instrumental error, and blank subtraction.
16 Optional post-processing incorporates common ion selection across replicate measurements
17 and shoulder ion removal. The scheme allows both positive and negative direct infusion
18 electrospray ionisation (ESI) and atmospheric pressure photoionisation (APPI) acquisition with
19 the same programs. An example application to atmospheric organic aerosol samples using an
20 Orbitrap mass spectrometer is reported for both ionisation techniques resulting in final spectra
21 with 0.8% and 8.4% of the peaks retained from the raw spectra for APPI positive and ESI
22 negative acquisition, respectively.



24 Highlights

- 25
- 26 • Ultra-high resolution mass spectrometry processing scheme from acquisition to analysis.
 - 27 • Method implementable with APPI and ESI ionisation in both polarities.
 - 28 • Example application for environmental samples showing >90% peak filtering.

29 Keywords

30 UHRMS, ESI, APPI, Environmental samples, direct infusion, Orbitrap

* Corresponding author: Markus Kalberer, markus.kalberer@atm.ch.cam.ac.uk

31 1 Introduction

32 Obtaining correct elemental composition of unknown molecules by mass spectrometry is a
33 challenge despite advances in instrumentation and data processing algorithms [1]. Ultra-high-
34 resolution mass spectrometry (UHRMS), coupled with soft ionisation techniques, most
35 commonly electrospray ionisation (ESI), can provide a detailed molecular composition for a
36 large, complex sample [2,3] being able to identify many distinct peaks at a given nominal mass.
37 Manual data processing and formula assignment is extremely time consuming [1,4] so
38 automatic algorithms have been developed that generally include noise elimination and blank
39 subtraction steps [5,6]. Noise filtering and blank subtraction is challenging for analysis in direct
40 infusion without prior chromatographic separation as the ion intensities may not be directly
41 related to the concentration of the molecules in the sample [7].

42 Constraints on allowed chemical elements and number of atoms are used when assigning
43 molecular formulae due to chemical reasons and computational limits. In a molecule containing
44 only carbon and hydrogen the “rule of 13” can be used to limit the number of carbon atoms, in
45 which the nominal mass is divided by 13 and the numerator gives the number of carbon atoms
46 and the remainder gives the number of hydrogens [8,9]. However, natural organic matter is
47 mainly composed of C, H, O and N with minor contributions from S and P, the latter being a
48 quantitatively non-significant component and often not considered [4,10,11]. The number of
49 possible solutions for an elemental formula increases largely if non-oxygen heteroatoms are
50 taken into consideration for calculation. Calculating unique elemental compositions is not
51 always possible [1] when acquiring data with high mass accuracy and resolution, especially as
52 increasing the molecular mass of analytes increases the number of possible molecular formula
53 assignments exponentially [4]. In addition, because mass spectrometry does not directly
54 provide structural information, these molecular formulae may represent any of several
55 structural isomers.

56 In order to automatically constrain the large number of possible candidates, rules have been
57 developed to select the most likely and chemically meaningful molecular formulae [1]. An
58 important constraint for restricting formulae to those that are likely to exist in nature is
59 including element ratios, especially the H/C ratio which, in most cases does not exceed $H/C > 3$
60 [1]. Similar restrictions can be put on the O/C (taking acidic polysaccharides as the upper limit
61 for molecular oxygen content $O/C < 1.3$) [4] and other heteroatoms to carbon ratios [11–13].
62 Additional constraints can be applied based on double bond equivalent (DBE), which indicates
63 the number of rings and double bonds in a molecule and is a measure of the degree of
64 unsaturation in a given compound [10]. Neutral molecules must have a DBE with an integer
65 value [8]. However, the elements N, S, and P may have different valences depending on their
66 chemical environment so constraints based on DBE values need to be used with caution [1].
67 Formulae are often filtered based on the “nitrogen rule” [1,4,14]. Neutral molecules containing
68 an odd total number of ^{14}N atoms always exhibit an odd nominal mass [4]. The nitrogen rule
69 derives from the fact that chemical elements with even nominal mass have an even valence,
70 while elements with odd mass have an odd valence, with the exception of nitrogen [3]. The
71 majority of data processing methods for Fourier transform ion cyclotron resonance technique
72 (FT-ICR) and OrbitrapTM mass spectrometers check for the presence of isotopes rather than
73 using isotopic ratios for formula assignments [4,10,15]. Other mass spectrometers, *e.g.*

74 TOF-type, often use isotope patterns for compound identification which tend to provide more
75 reliable assignments compared to FT-ICR and Orbitrap™ mass spectrometers [16].
76 Once chemically meaningful formulae have been filtered, more than one possible formula per
77 peak may still exist, especially at high m/z . In order to select the most meaningful formula
78 assignment, two general strategies have been applied: a “best-fit” approach, in which the
79 formula with the closest match between theoretical mass and observed mass is selected, and a
80 “formula extension” approach, in which chemical and structural relationships among
81 compounds are taken into account for formula assignment, *e.g.* by looking for homologous
82 series based on Kendrick mass defects [5,14,17,18]. In the first case, possible incorrect
83 assignments may arise from inaccuracies in the measured masses [19–21]. In the second case,
84 it has previously been observed that, for example, atmospheric oxidation reactions involving
85 S- and N-containing functional groups may lead to a wide variety of products that do not
86 produce homologous series, risking the removal of potentially correct assignments [22].
87 Most of the methods and currently available algorithms were developed based on ESI and
88 therefore they rely on the assumption that ionisation is accompanied by protonation,
89 deprotonation or adduct formation [20,23,24]. Other ionisation techniques, such as atmospheric
90 pressure chemical ionisation (APCI) and atmospheric pressure photoionisation (APPI) are
91 becoming increasingly common for less polar and apolar organic compounds [20,25]. In APPI,
92 detection of molecular ions (as radical cations or anions) over quasi-molecular ions is common
93 [20,25] so there is a need to develop new algorithms that take into account the formation of
94 molecular ions.
95 Here we developed a code to filter molecular formula assignments that i) can be applied to
96 different soft-ionisation techniques like ESI, and APPI in both positive and negative ionisation,
97 ii) takes into account formation of molecular ions, quasi-molecular ions and Na adducts, iii)
98 uses a novel method for mass shift and noise estimation and iv) can be used with two different
99 blank subtraction methods. Many steps of the scheme are widely used in mass spectrometry
100 studies, but direct comparisons between methods are difficult as detailed procedures are often
101 not available in the literature. Aspects of the approach described here have been previously
102 applied in studies of environmental [12] and biological samples [26].

103 **2 Pre-processing**

104 The following discussion is based on the use of an Orbitrap™ mass spectrometer (LTQ
105 Orbitrap Velos, Thermo Scientific™, Bremen, Germany) with the proprietary software
106 Xcalibur™ 2.1-3.0 (Thermo Scientific™, Bremen, Germany) henceforth referred to simply as
107 *Xcalibur*. The steps taken, however, apply for the general processing of mass spectra with any
108 spectrometer/software combination. While the choice of mass spectrometer does not affect the
109 presented procedure, the resolution of the spectrometer will influence the accuracy of the final
110 spectra with higher resolutions providing clearer peak separation and higher confidence in the
111 molecular formula assignment [27] which is particularly important for complex ambient
112 measurements [28].

113 2.1 Data acquisition

114 Ion transmission and ion collection efficiencies in an Orbitrap™ mass spectrometer strongly
115 depend on the m/z scan range [29]. Therefore, to avoid the loss of the ions at the high or low
116 end of the selected mass range (*e.g.* m/z 50-1000), it can be split into several overlapping scan
117 ranges [29]. Each of these scan ranges are independently processed and subsequently
118 recombined. Both sample (*i.e.* the spectrum of interest) and blank (*i.e.* a reference spectrum)
119 spectra are acquired under the same conditions. The type of blank will vary based on
120 application but can include solvent, procedural, and field blanks. The blanks are later used to
121 remove contaminants (Section 3.1.3).

122 The reproducibility of the peak centroids (the mode of the intensity distribution of an individual
123 peak) and their magnitudes is an important aspect of the direct infusion ESI-UHRMS method,
124 especially for low intensity ions due to competitive ionisation and matrix artefacts. Therefore,
125 besides applying noise threshold corrections (Section 3.1.1), it is also important to consider
126 instrumental replicate measurements. For example, natural organic matter sample replicates
127 are considered reproducible with ESI FT-ICR-MS if a minimum of 67% of threshold-corrected
128 peaks are common among replicates [30]. Additionally, longer acquisition times are desirable
129 as the signal-to-noise ratio improves in proportion to the square root of time [31]. The following
130 discussion therefore assumes multiple replicates (used to filter out uncommon assignments as
131 discussed in Section 3.1.6) and multiple scans to provide a reasonable average (as discussed in
132 Section 2.2).

133

134 2.2 Data pre-treatment

135 Each acquired mass spectra is averaged into one spectrum to reduce the noise level. Molecular
136 assignments are performed using *Xcalibur* software applying constraints on the allowable
137 number of each element, the maximum number of possible formulae to assign, and the
138 maximum mass error.

139 The restriction on the number of elements varies based on application (see Table 1) but the
140 current procedure permits limits on ^{12}C , ^{13}C , ^1H , ^{14}N , ^{16}O , ^{32}S , and ^{34}S for both ionisation
141 polarities. For positive ionisation modes up to one ^{23}Na atom is additionally allowed. The
142 absolute element limits are generally determined by dividing the mass range through the
143 element mass; by using the developed set of formulae that were derived from the National
144 Institute of Standards and Technology (NIST), Wiley mass spectra, and the Dictionary of
145 Natural Products (DNP) entries as discussed in Kind and Fiehn [1]; and considering the
146 presence of oligomeric compounds. For the heavy isotopes of carbon and sulphur, their natural
147 abundances are so low that the relative abundance of molecules with multiple heavy isotopes
148 is low (*e.g.* the natural abundance of molecules with more than two ^{13}C atoms is below
149 detection limits [32]). An *a priori* knowledge of expected elements is important as excluding
150 elements potentially removes correct formula assignments and including additional elements
151 risks calculating unrealistic, but lower mass error, assignments [4].

152 The first 5 (or more) mathematically possible elemental formulae (depending on mass range
153 and instrument accuracy) with the lowest mass error value within a given mass tolerance (up
154 to ± 6 ppm) are exported. A wide mass tolerance is used to cover for the observed non-

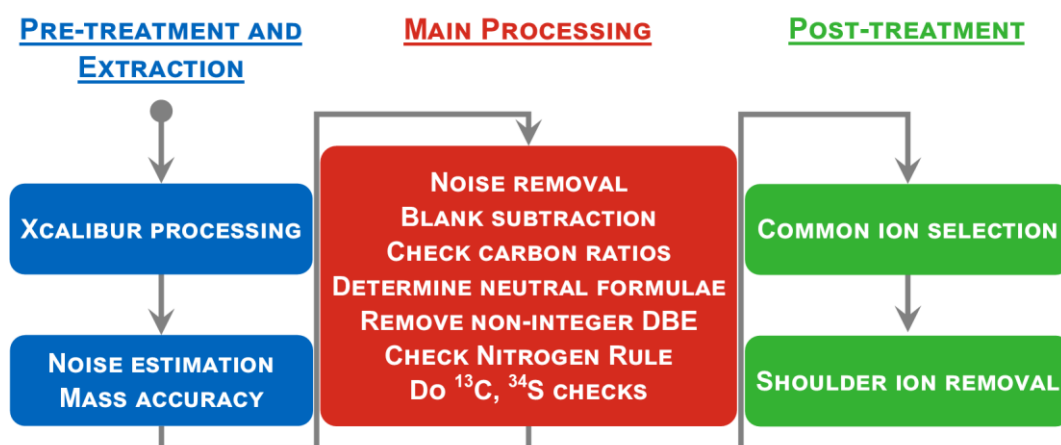
155 systematic mass shift, which seems to be dependent on several factors *e.g.* the sample matrix,
156 the ion intensity of individual ions within this matrix, and the mass range. At later data
157 treatment stages the mass shift is estimated (Section 3.1.2) and subsequently corrected during
158 the main processing to significantly narrow down the mass tolerance. Lock masses can be used
159 during acquisition to reduce (but not remove) the observed mass shift [19,31] to improve
160 formula assignment [7,33].
161

162 3 Processing and discussion

163 After pre-treatment *via Xcalibur* each spectrum requires a number of steps to remove unlikely
164 formula assignments associated with, for example, instrument noise and sample contamination.
165 These checks, amongst others, are included in the data treatment discussed below.

166 3.1 Data treatment

167 The overall data treatment procedure can be split into three major stages: pre-treatment and
168 extraction, main processing, and post-treatment. Pre-treatment and extraction consists of the
169 steps using *Xcalibur*, discussed above, as well as the initial extraction of mass shift and noise
170 estimations for both sample and blank spectra. The main processing stage includes all major
171 filter and blank subtraction processes. Finally, post-treatment consists of common ion selection
172 (over several repeated measurements of the samples) and shoulder ion removal. The stages and
173 major steps within each stage are listed in Figure 1. The following sections will describe each
174 step in more detail. Individual rounded rectangles denote separate processing scripts which are
175 predominantly written in Mathematica 10.4 (Wolfram Research Inc., UK), henceforth referred
176 to as *Mathematica*, except for the *Xcalibur* processing step which is manually processed.
177
178



179
180 **Figure 1** Schematic overview of the three main processing stages of data processing: pre-treatment and
181 extraction (left), main processing (centre), and post-treatment (right). Each rounded rectangle denotes a
182 separate script being used within each stage. The order of steps within the “Main Processing” stage varies
183 slightly based on blank subtraction method (see Section 3.1.3). DBE ≡ double bond equivalent.

184 3.1.1 Evaluation of noise level

185 Previous studies have used an extremely wide variety of noise levels evaluated from signal-
186 free regions in the mass spectra [34]. Hawkes *et al.* [35] observed that intensity of noise
187 increases with m/z in FT-ICR but it is constant in the Orbitrap™ between 150-2000 Da. Signal-
188 to-noise ratios (S/N) ranging between 2 and 20 [2,7,8,10,13,21,24,34,36] have been used which
189 highlights the challenge of reliable noise level estimation in a wide mass range. Other studies
190 simply removed the lowest 10% of peaks (based on intensity) [6] while a S/N of 4 effectively
191 removed peaks with relative intensity below 0.5% in a previous study [24].

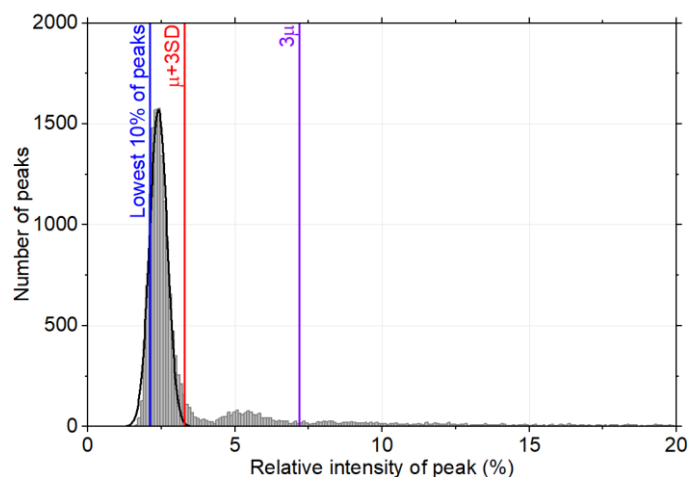
192 The choice of S/N or a cut-off based on intensity seems arbitrary since intensity of peaks is
193 affected by competitive ionisation and cannot be used to infer concentration of compounds
194 when using direct infusion [7,36].

195 Sleighter *et al.* [30] compared reproducibility in peak detection at different S/N thresholds.
196 They showed that reproducibility increases with increasing S/N cut-off from 3 to 10. However,
197 they point out also that using a strict S/N threshold is not adequate for establishing peak
198 detection reproducibility because well-defined peaks could go undetected just below the
199 defined threshold. They suggest using a lower S/N of 2.5 when looking for common peaks in
200 the other replicates.

201 In our method, for each processed spectrum (sample and blank), the noise level is estimated
202 based on fitting a normal distribution to a histogram of intensities. The process is visualised in
203 Figure 2 which shows a typical intensity distribution. Histogram bin sizes are selected based
204 on the Freedman-Diaconis rule [37] to ensure the histogram is representative without excessive
205 computer processing. The noise intensity is characterised by a bi-modal normal distribution
206 (Figure 2) which has also been seen by Zhurov *et al.* when describing their alternative approach
207 to the “N sigma” methodology for determining noise levels [38]. The first mode corresponds
208 to the lowest intensity peaks in the MS probably associated with thermal noise [31]. The second
209 mode may correspond to a higher intensity chemical noise [6]. For example, we observe several
210 shoulder ions present nearby to high intensity peaks in the mass spectrum. These artefact ions
211 may have intensities similar to analyte peaks with low concentrations or ionisation efficiency
212 which make up the second mode of the histogram. It is therefore difficult to discriminate
213 between high intensity noise and low intensity analyte peaks.

214 For this reason, the noise level, which is subsequently used to remove peaks during the main
215 processing stage, is estimated as the mean plus three standard deviations based on the fit of the
216 first mode whenever it is possible to acquire at least three instrumental repetitions. This
217 approach follows the “N sigma” methodology which implements a noise level equal to the
218 mean plus N times the standard deviation (*i.e.* σ). Using $N = 3$ is conservative considering
219 typical N values are 3, 5, 6, or 8 [38]. The second mode is addressed during a later processing
220 stage where only peaks common in all (or a user-defined fraction) of the replicates are kept
221 (see Section 3.1.6) and these are subsequently filtered to remove shoulder ions at masses close
222 to high intensity peaks (see Section 3.1.7).

223 In contrast, when it is not possible to acquire three instrumental repetitions per sample, *e.g.* in
224 Liquid Extraction Surface Analysis (LESA) when analysing chemically heterogeneous
225 surfaces [26,39], the noise level is estimated as (at least) three times the mean in order to delete
226 all peaks appertaining to the second mode of the intensity distributions (Figure 2) assuming
227 they are high intensity chemical noise.



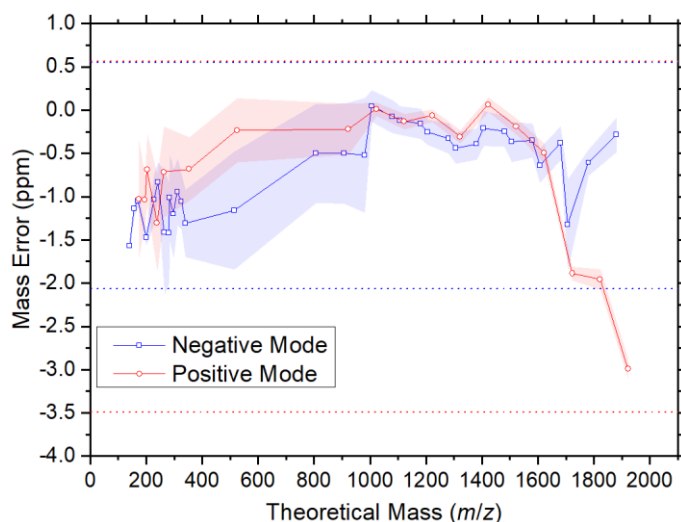
228
 229 **Figure 2** Sample intensity histogram used for noise level estimation using raw data (grey) and fitted curve
 230 (black). The noise level (red line) is typically set to the mean plus three standard deviations based on the
 231 fit. Under some conditions, when there are not enough replicates for common ion comparisons, the noise
 232 level may be set to three times the mean (purple). An additional noise level (blue line) marks the intensity
 233 level below which 10% of all peaks reside based on Wong *et al.* (2009) [6]. The intensities (along the *x*-axis)
 234 are given as a percentage of the maximum intensity peak in the spectrum.

235 3.1.2 Evaluation of mass shift

236 The mass shift of each mass spectrum is evaluated based on the measured mass errors of a
 237 minimum of ten known reference compounds from *a priori* knowledge (based on the source
 238 and polarity used) which are expected in either the sample or solvent. These reference peaks
 239 provide an estimate of the overall mass shift of a given spectrum (similar to the method
 240 described by Sleighter *et al.* [40]). The reference compounds should be selected to cover the
 241 mass range of interest due to the mass dependence of the mass error, as shown using calibration
 242 standards in Figure 3, and to be well resolved (or with a higher signal compared to adjacent
 243 peaks) in order to avoid inaccuracies in the mass shift evaluation. Mass errors are also
 244 dependent on the peak intensity [41] and the matrix but these factors are more difficult to
 245 account for and are not considered here systematically.

246 The errors, as originally calculated by *Xcalibur* in units of ppm, are summarised by their
 247 arithmetic mean, standard deviation, maximum, and minimum values to be used in subsequent
 248 processing steps. The mass errors are tested for outliers, to avoid skewed data, using Grubbs'
 249 Test [42] with a default confidence level of 99%. Any detected outliers are removed from the
 250 summary data. While calculating mass shifts, if several peaks are assigned the same chemical
 251 formula the highest intensity peak is assumed to have the correct assignment in order to avoid
 252 selecting shoulder ions (see Section 3.1.7). Functionally, mass shifts are used as a filter for
 253 removing assignments with errors outside an acceptable range (based on the
 254 minimum/maximum mass shift of known compounds) effectively restricting the initial wide
 255 range set in *Xcalibur*. In order to account for the intra-spectrum mass shift variability
 256 mentioned above, a conservative approach was implemented where the maximum and
 257 minimum mass shifts allowed for formula assignments can be increased (at the user's
 258 discretion) beyond the range observed for the reference compounds (see Figure 3) based on
 259 user input with typical values being on the order of 0.5 ppm. The mass error in the example
 260 shown in Figure 3 has a range of ca. 2.5 ppm for negative mode (from ca. 0.5 ppm to
 261 ca. -2 ppm) and about 4 ppm for positive mode (from ca. 0.5 ppm to ca. -3.5 ppm) ionisation

262 and illustrates that for most spectra the error is not symmetrically distributed around 0 ppm. At
263 this stage, mass shifts are only used to restrict the range of possible formula assignments
264 without explicitly selecting the correct formula assignment. Mass shifts may also be used
265 during blank subtraction to effectively realign the sample and blank peaks for comparison
266 purposes using the mean and standard deviation as discussed in the following section.
267



268 **Figure 3** Mass shift of known calibration mixture for ESI positive (red) and negative (blue). Markers denote
269 the mean value from at least three repeats with the shaded area showing the associated standard deviation.
270 Non-uniform, non-linear relationship is shown for both modes with mass errors peaking at extreme m/z
271 values. A minimum and maximum approach is used to account for this known relationship as shown with
272 the dashed limit (with an additional offset of 0.5 ppm included). Calibration standards used were the Pierce
273 ESI Negative Ion Calibration Solution (Thermo Scientific) and the Pierce LTQ Velos ESI Positive Ion
274 Calibration Solution (Thermo Scientific). Data taken from three m/z scan ranges: 80-600, 150-1000, and
275 200-2000.
276

277 3.1.3 Blank subtraction

278 Blank subtraction aims to remove any peaks from the sample spectra that are also present in
279 the blank spectra. Two approaches have been studied based on either the final assigned
280 chemical formulae or the corrected masses of each peak. The first method, based on chemical
281 formulae, processes the sample and blank independently and removes matching assigned
282 molecular formulae that have a sample-to-blank ratio below a user-defined limit (*e.g.* 10). The
283 general approach of processing the sample and blank independently before comparison is
284 common in literature [5,43–45], although the removal criteria vary.

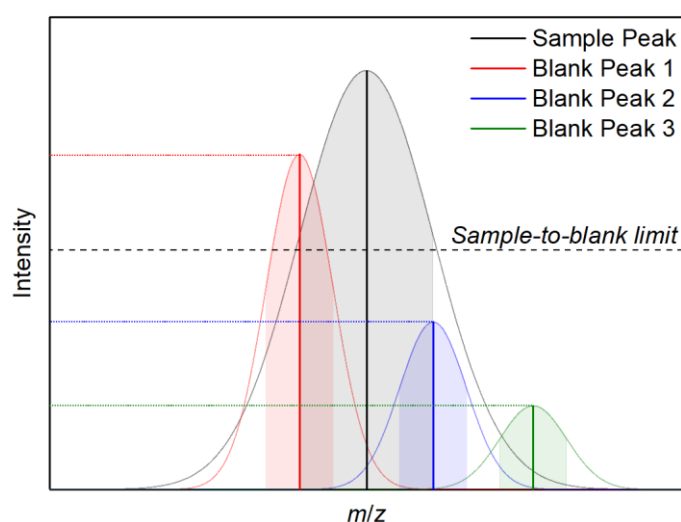
285 The second method checks whether a sample and blank peak are statistically equivalent and
286 removes the sample peak if its intensity is not larger than the sample-to-blank ratio. In this
287 second method the first check, based on a two sample *t*-test, determines whether the sample
288 and blank are within one standard deviation of one another (based on the mass shift standard
289 deviations and correcting for the mean mass shifts to account for inter-spectra variability). If
290 this is true, as demonstrated by blank peaks 1 and 2 in Figure 4, then the second check is
291 performed. The second check determines if the ratio of sample peak intensity to blank peak
292 intensity is below a certain level, *i.e.* the blank is above a certain sample-to-blank limit as
293 illustrated in Figure 4 by a dashed line. If this is true, as in the case of the example blank peak
294 1, then the sample peak is removed from the dataset. Otherwise, the sample peak is retained

295 under the assumption that the blank peak is not the same compound (*e.g.* blank peak 3 in Figure
296 4) or that there is significantly more of the given compound in the sample as compared to the
297 blank (*e.g.* blank peak 2 in Figure 4).

298 Although both methods produce similar results there are minor differences between the two.
299 The formula-based approach has a strong dependence on the number of assigned formulae
300 during the initial *Xcalibur* processing that may lead to false positives if the same formula is not
301 assigned within the blank spectrum due to the influence of different mass errors on the *Xcalibur*
302 assignment algorithm. That is, for a given number of formula assignments, a small change in
303 measured mass may result in different sets of formulae for the sample and blank. After
304 processing this may allow different final assignments for effectively the same peak resulting in
305 the sample peak being incorrectly kept. On the other hand, the mass-based approach may be
306 limited by the mass error variation. When the standard deviation of the mass shift is high (*i.e.*
307 greater than instrument accuracy), as determined in Section 3.1.2, false negatives due to
308 overzealous blank subtraction may occur. This is due to a wider range of blank peaks being
309 compared to the sample and potentially satisfying the conditions for removal. The result of the
310 difference is typically <5% for the two blank subtraction methods with the formula-based
311 approach having additional false positive assignments with reasonable mass shift variability.
312 Such a comparison, however, is strongly dependent on the number of formulae assigned per
313 peak and the mass shift standard deviation.

314 The sample-to-blank ratio, when used, is largely arbitrarily selected. Previous studies have used
315 ratios of 10 [2,7,43], up to effectively ∞ [30,45] (*i.e.* everything in the blank is removed from
316 the sample). Rincon *et al.* [44] had a hybrid approach where peaks below a sample-to-blank
317 ratio of 1 would be removed (if within 2 ppm), otherwise the blank intensity was subtracted
318 from the sample assuming the matrix effects were similar for both the sample and blank. For
319 this reason, the sample-to-blank ratio is an adjustable user input during processing.

320



321
322 **Figure 4 Schematic of the blank subtraction process for a sample peak (black) being compared to three**
323 **nearby blank peaks (red, green, blue). The first check determines whether the peaks overlap within one**
324 **standard deviation (*i.e.* shaded regions) which is the case for blank peaks 1 and 2. The second check is**
325 **whether the sample-to-blank ratio is above a specified sample-to-blank limit which is only true for blank**
326 **peak 1. Therefore, the sample peak is removed due to the presence of blank peak 1.**

327 *3.1.4 Additional exclusion criteria*

328 In addition to noise removal and blank subtraction there are supplementary filters that the main
329 processing code performs.

330

331 **Carbon ratios**

332 Previous studies (Table 1) have used carbon ratios to eliminate compounds unlikely to naturally
333 exist in the sampled environment. As such, the code allows for control of the exclusion of
334 certain O/C, H/C, N/C, and S/C ratios based on user input. Phosphorous is not currently
335 considered due to its unlikelihood for being a significant component of atmospheric samples
336 (for which this code was initially developed) [10] but is more important for water [46] and soil
337 [11] samples. These limits have seen variability between different references as shown in Table
338 1 in both value and the choice of ratios used for filtering as expected for varying environments.
339 In the current processing scheme, we refer to the H/C ratio of the neutral molecular formula
340 which is calculated differently depending on ionisation source type (ESI vs. APPI) and polarity
341 (positive vs. negative). We assume that in ESI the dominant ions are $[M+H]^+$ and $[M+Na]^+$ in
342 positive ionisation and $[M-H]^-$ in negative ionisation. Conversely, in APPI molecular ions $[M]^+$
343 and $[M]^-$ are also present [25] in addition to quasi-molecular ions and sodium adducts. In order
344 to distinguish between molecular ions ($[M]^+$, $[M]^-$) and quasi-molecular ($[M+H]^+$, $[M-H]^-$) ions
345 we use DBE values. In APPI, when the DBE of the ion is a non-integer we assume it is a quasi-
346 molecular ion while when the DBE is an integer we assume it is a molecular ion.

Table 1 Allowable atom count and carbon ratio ranges from literature to filter out non-naturally existing chemical formula assignments.

Reference	Instrument	Allowable atom count										Allowable ratio range																
		¹² C		¹ H		¹⁴ N		¹⁶ O		²³ Na		³² S		³⁴ S		³¹ P		¹³ C		O/C		H/C		N/C	S/C	P/C	(S+P)/C	
		≥	≤	≥	≤	≥	≤	≥	≤	≥	≤	≥	≤	≥	≤	≥	≤	≥	≤	≥	≤	≥	≤	≤	≤	≤	≤	
Koch <i>et al.</i> 2005 [8]	FT-ICR	-	100	-	200 ¹	-	-	-	50	-	1	-	-	-	-	-	-	0	1.5	0.3	-	-	-	-	-	-	-	
Koch and Dittmar 2006 [47]	FT-ICR	1	100	1	200	0	10	0	50	0	1	-	-	-	-	-	-	0	1.2	0	2.2	0.5	-	-	-	-	-	
Kind and Fiehn 2007 [1] ²	FT-ICR	Allowable atom count varies based on mass range and reference library															0	1.2	0.2	3.1	1.3	0.8	0.3	-				
Koch <i>et al.</i> 2007 [4] ³	FT-ICR	0	∞	0	∞	0	30	0	∞	-	-	0	2	-	-	0	2	-	1	0.3 ⁴	-	1	-	-	-	-	-	
Wozniak <i>et al.</i> 2008 [10] ⁵	FT-ICR	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	0	1.2	0.3	2.25	0.5	0.2	0.1	0.2	-	-	
Altieri <i>et al.</i> 2009 [48]	FT-ICR	-	∞	-	∞	-	15	-	15	-	-	-	1	-	-	-	1	-	5 ⁴	0.3	-	2 ⁴	-	-	-	-	-	
Bateman <i>et al.</i> 2009 [45]	LTQ-Orbitrap	-	50 ⁴	-	100 ⁴	-	-	-	60 ⁴	0	1	-	-	-	-	-	-	0.05	1.3	0.7	2	-	-	-	-	-	-	
Schmitt-Kopplin <i>et al.</i> 2010 [27]	FT-ICR	-	20	-	30	-	5	-	6	-	-	-	1	-	-	-	-	0	1	0	2n+2	-	-	-	-	-	-	
Stubbins <i>et al.</i> 2010 [49]	FT-ICR	-	50	2	2c+2	-	-	0	c+2	-	-	-	-	-	-	-	-	0	1.2	0.333	2.25	0.5	0.2	0.1	0.2	-	-	
Fuller <i>et al.</i> 2012 [39]	LTQ-Orbitrap	1	20	-	-	-	-	-	-	0	1	0	1	-	-	-	-	0	3	0.2	3	1	-	-	-	-	-	-
Rincón <i>et al.</i> 2012 [44]	LTQ-Orbitrap	-	35	-	75	-	7	-	25	-	1	-	7	-	-	-	-	0	5	0.3	7	6	-	-	-	-	-	-
Kourtchev <i>et al.</i> 2013 [12]	LTQ-Orbitrap	-	100	-	200	-	5	-	50	-	-	-	2	-	1	-	-	0	1.5	0.3	2.5	0.5	0.2	-	-	-	-	-
Ohno and Ohno 2013 [11]	FT-ICR	8	50	8	100	0	5	1	30	-	-	0	3	-	-	0	2	-	1.2	0.3 ⁴	2.25 ⁴	0.5 ⁴	0.2 ⁴	0.1 ⁴	0.2	-	-	
Kourtchev <i>et al.</i> 2014 [7]	LTQ-Orbitrap	-	100	-	200	-	5	-	50	-	-	-	2	-	1	-	-	0	1.2	0.3	2.5	0.5	0.2	-	-	-	-	-
Fooshee <i>et al.</i> 2015 [50]	LTQ-Orbitrap	1	80	2	140	-	-	0	50	0	1	-	-	-	-	-	-	0	1.2	0.5	2.2	-	-	-	-	-	-	-
Lu <i>et al.</i> 2015 [13]	FT-ICR	1	50	2	100	0	6	0	30	-	-	0	2	-	-	-	-	0	1.2	0.35	2.25	0.5	0.2	-	-	-	-	-
Herzprung <i>et al.</i> 2016 [21]	FT-ICR	-	100	-	-	-	5	-	80	-	-	-	-	-	-	-	-	0	1	0.3	2n+2	-	-	-	-	-	-	-
Wang <i>et al.</i> 2016 [43]	Q-Exactive	1	40	2	80	0	3	0	40	-	-	0	2	-	-	-	-	0	3	0.3	3	0.5	0.2	-	-	-	-	-

Lower-bound for ratios is assumed to be zero unless otherwise specified.

Unspecified values denoted by dash (-).

c, *h*, and *n* denote the number of carbon, hydrogen, and nitrogen atoms, respectively.

Inclusion of sodium is for positive ionisation only.

¹ Also includes an additional filtering of $h < 2c + 2$.

² Kind and Fiehn 2007 [1] ratio values are covering the 99.7th percentile.

³ Koch *et al.* 2007 [4] atom count ranges are varied. Quoted values are the most inclusive set.

⁴ Quoted value is reported as an exclusive range (*e.g.* < rather than ≤).

⁵ Wozniak *et al.* 2008 [10] removed all phosphorus containing compounds after initial filtering.

348
349
350
351
352
353
354
355
356

357 Double bond equivalent values

358 Neutral formulae with non-integer (*e.g.* charged molecule [8] and radicals [1]) or negative
359 double bond equivalent (DBE) values are removed from the final peak list. The DBE value is
360 a metric for unsaturation of a given compound based on the number of rings and double bonds
361 [10]. DBEs are calculated using the following:

$$362 \text{DBE} = 1 - \frac{h}{2} + \frac{n}{2} + \frac{s}{2} + c$$

363 where *c*, *h*, *n*, and *s* correspond to the number of atoms (*i.e.* C, H, N, and S) in any given
364 chemical assignment $C_cH_hN_nO_oS_s$ [10,12]. Similar to the carbon ratios, the DBE calculation
365 varies based on source and polarity because of the assigned hydrogen count.

367 Nitrogen rule

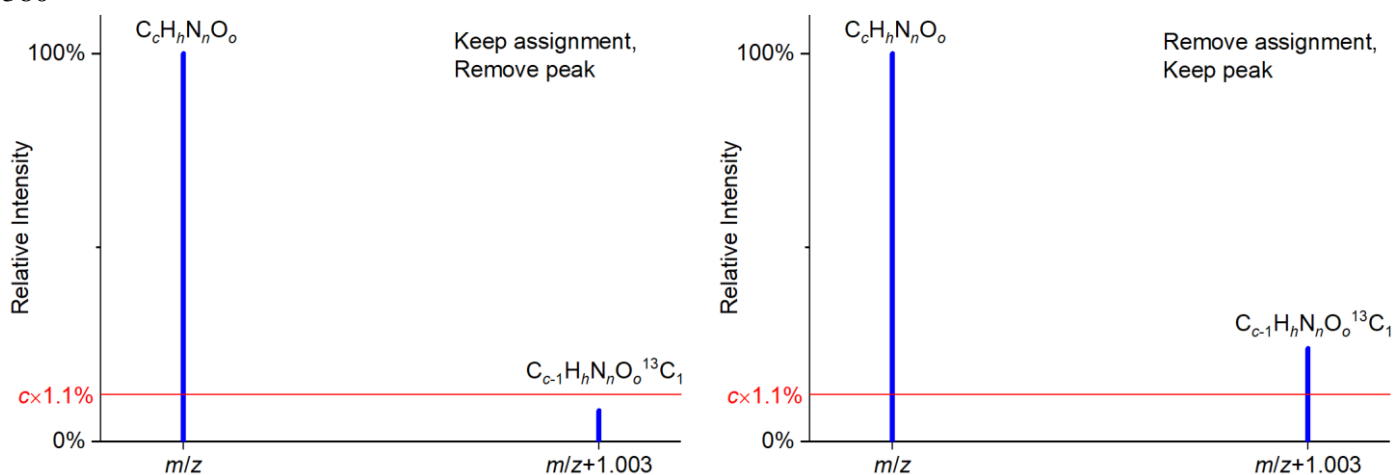
368 Any peaks failing the Nitrogen rule [8], which states that a neutral compound with an odd
369 nominal mass has an odd number of nitrogen atoms [1], are removed.

370

371 Carbon and Sulphur isotopes

372 When assignments contain ^{13}C and/or ^{34}S the filtering process checks for a peak with a
373 matching composition containing only the lighter isotope. If the intensity ratio of
374 heavier-to-lighter isotope was greater than a user-defined factor (*e.g.* 1.2) of the natural isotopic
375 abundance (1.1% and 4.5% for ^{13}C and ^{34}S , respectively [51]) then the isotopic assignment is
376 removed (*i.e.* the assignment with the next smallest mass delta was used). Otherwise, the
377 isotopic assignment is considered the correct assignment and the entire peak is removed as it
378 is chemically equal to the peak containing only the lighter isotope. This process is illustrated
379 in Figure 5.

380



381 **Figure 5** Isotopic removal schematic with $^{13}\text{C}/^{12}\text{C}$ example showing peaks with ^{12}C isotopes only (left) and
382 with one ^{13}C atom (right) for an arbitrary assignment $C_cH_hN_nO_o$ using a ^{13}C natural abundance of 1.1%.
383 When the relative intensity of the isotope is below (left panel) the $c \times 1.1\%$ limit, denoted by the red line, the
384 isotope assignment is considered correct and the peak is removed as it is redundant. Otherwise (right
385 panel), the assignment is deemed incorrect and removed while the peak remains using the assignment with
386 the next smallest mass error.

387 This approach takes into account the possibility that more than one compound may contribute
388 to the same observed peak within instrumental accuracy and this may be the case when the
389 intensity of the isotopic peak is higher than what expected from its natural abundance.
390 Conversely, the method described by Wozniak *et al.* [10] removed any peak 1.003 m/z units
391 above another peak under the assumption that the peak at the higher m/z is always the ^{13}C
392 isotope. Ohno *et al.* [15] removes peaks with intensity lower than 50% of the lighter isotope
393 ion. Heavier-to-lighter isotope ratios tend to be underestimated by both FT-ICR [4] and
394 Orbitrap™ [52] analysers compared to theoretical ratios. Given the likelihood of
395 underestimation, the aforementioned isotopic ratio factor can safely be set to 1 (*i.e.* use the
396 natural abundance) but varying the value can be shown to still influence the number of final
397 peaks – especially if the abundance is lowered (see Figure S1 in Supplementary Material).

398 3.1.5 Duplicate removals

399 At this stage of the processing there may still be multiple assignments for a single peak within
400 the derived experimental mass error range (as described in Section 3.1.2). An option within the
401 *Mathematica* script allows for duplicate removal where the assignment with the smallest
402 absolute mass error, after mass shift correction, is kept as the true assignment. The option can
403 be selected at the user’s discretion depending on the scenario. In general, if the resolution and
404 accuracy of the instrument used do not allow the identification of a unique formula assignment
405 for a given peak in the mass spectrum, different approaches may be considered. Those include:
406 (i) selecting the assignment with the smallest mass error (option available within the script);
407 (ii) selecting the assignment using a “formula extension” approach based on Kendrick mass
408 defects (not implemented within the script); (iii) keeping all possible assignments (option
409 available within the script if the duplicate removal option is not used); and (iv) removing all
410 peaks for which multiple assignments are still present (possible to do manually after the data
411 processing). Although the first approach is implemented in the processing scheme, care should
412 be taken since the best formula assignment may not always be assignment with the lowest mass
413 error [53]. The “formula extension” approach based on Kendrick mass defects was not
414 implemented because, as mentioned in the introduction, S- and N-containing functional groups
415 do not necessarily form a homologous series in a Kendrick mass defect plot [22].

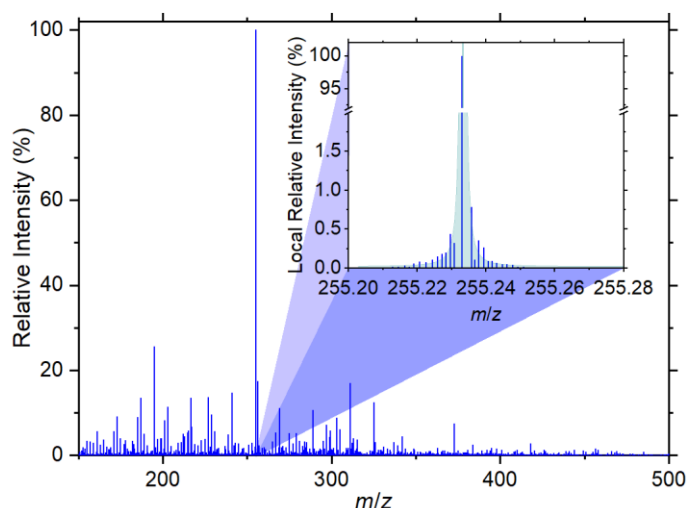
416 3.1.6 Common ion selection

417 After several replicate mass spectra (minimum of 3) of each sample are processed through the
418 main processing stage of Figure 1, they are simultaneously compared for common ions. This
419 common ion selection process filters out any peaks that do not exist in all (or some fraction) of
420 the processed replicates as chosen by the user. The output intensity of the common ions is the
421 average of the replicates. Common ion selection removes the second mode visible in the noise
422 histograms (Figure 2), as discussed above, as those peaks are predominantly noise.

423 3.1.7 Shoulder peak removal

424 Shoulder peaks are artefacts of the mass spectrometer’s processing produced during the Fourier
425 transform calculation [54]. Figure 6 shows a scenario of shouldering (inset) which is
426 highlighted by a shaded Lorentzian curve fit. The high intensity peak is bordered by several
427 low intensity shoulder peaks. Given the difficulty to identify these artefacts, a conservative
428 approach is used to remove apparent shoulders based on the assumption that shoulders are more

429 likely as peak intensity increases. If a peak is intense enough (*e.g.* >1,000,000), any
430 neighbouring local ions (*e.g.* within ± 0.01 m/z) that are less than a specified percentage (*e.g.*
431 1%) of the local major peak are considered shoulders and removed. More intense peaks are
432 kept as they could still be considered true peaks. The local peak intensity, mass range, and
433 shoulder percentages are adjustable by the user.
434



435
436 **Figure 6** An example of shoulder peaks surrounding a high intensity peak (inset). A shaded Lorentzian fit
437 highlights the shoulder peaks that are all below 1% relative intensity and within ± 0.02 m/z . The example
438 shown was analysed in the ESI negative ionisation mode.

439 **3.2 Sample Application: UHRMS analysis of the organic fraction of urban PM_{2.5}** 440 **samples**

441 The sample application of the data processing procedure is based on a PM_{2.5} filter taken on 30th
442 May 2014 over 24 hrs at an urban background site in the city centre of Padua (Italy), located
443 in the polluted Po Valley. More details of the sampling site and procedures for sample
444 collection are reported elsewhere [55]. The filter was extracted in methanol using the procedure
445 already described elsewhere [7] and analysed in both APPI and ESI in both polarities. Here we
446 show the sample processing performed for APPI positive and ESI negative ionisation modes.

447 *3.2.1 Instrumental analysis/Data acquisition*

448 Samples were analysed with a high resolution LTQ Orbitrap Velos mass spectrometer (Thermo
449 Scientific™, Bremen, Germany) equipped with a TriVersa Nanomate® chip-based ESI source
450 (Advion Biosciences, Ithaca NY, USA) and APPI Ion Max source (Thermo Scientific™,
451 Bremen, Germany) with a Syagen Krypton lamp emitting photons at 10.0 eV and 10.6 eV.

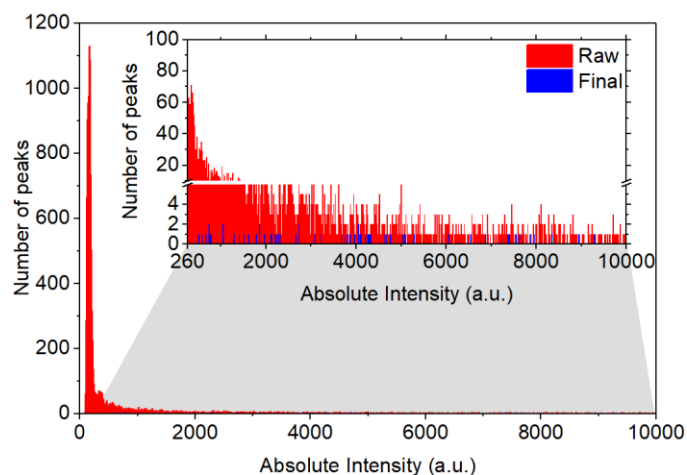
452 The direct infusion negative nanoESI parameters were as follows: ionization voltage -1.6 kV,
453 back pressure 0.8 psi, capillary temperature 275 °C, S-lens RF level 60%, sample volume 8 μ L.
454 For analysis in APPI, methanolic extracts doped with 10% toluene were infused at a flow rate
455 of 10 μ L/min, with a source temperature of 200 °C, a sheath gas flow of 0 L/min, an auxiliary
456 gas flow of 5 L/min, and a sweep gas flow of 10 L/min. The mass analyser was calibrated
457 before the analysis using Pierce LTQ Velos ESI Positive Ion Calibration Solution (Thermo
458 Scientific) and Pierce ESI Negative Ion Calibration Solution (Thermo Scientific). The mass
459 accuracy of the instrument was checked before the analysis and was below 0.5 ppm. The
460 instrument mass resolution was set at 100,000 (at m/z 400). Each sample was analysed in the

461 m/z ranges 100–650 and 150-900, acquiring four repeats for 60 seconds each (~40 scans) in
462 centroid mode.

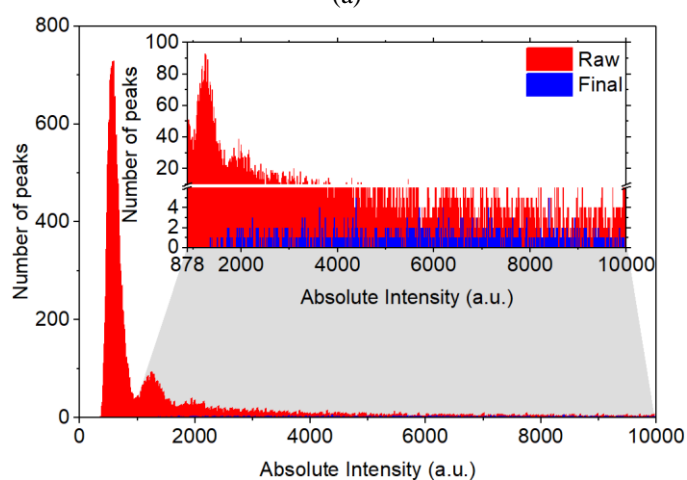
463 The four sample repeats processed for the m/z range of 150-900 are discussed here for each
464 ionisation mode. Within *Xcalibur* the chemical assignments were determined for up to 10
465 formulae per peak allowing a maximum mass error of ± 6 ppm. The formula assignment was
466 performed using $1 \leq {}^{12}\text{C} \leq 75$, $0 \leq {}^{13}\text{C} \leq 1$, $1 \leq {}^1\text{H} \leq 180$, $0 \leq {}^{16}\text{O} \leq 50$, $0 \leq {}^{14}\text{N} \leq 30$,
467 $0 \leq {}^{32}\text{S} \leq 2$, $0 \leq {}^{34}\text{S} \leq 1$. The positive ionisation mode additionally allowed for up to one sodium
468 atom.

469 3.2.2 *Data processing*

470 The data from each ionisation mode was processed using a noise level based on the mean plus
471 three standard deviations definition (as discussed in Section 3.1.1). The resulting noise levels
472 were calculated to be 260 and 878 for the APPI positive and ESI negative modes, respectively.
473 A comparison of the intensity histograms before and after the entire processing are shown in
474 Figure 7 as a means of examining the change in high count, low intensity peaks typically
475 associated with noise. As expected, the secondary mode was removed during the common ion
476 stage of post-processing and the entire histogram was reduced to sub-20 counts for all
477 intensities.



(a)



(b)

478 **Figure 7 Intensity histograms before (red) and after (blue) processing for (a) APPI+ and (b) ESI- samples.**
 479 **Inset zooms to the range above the preset ‘noise limit’ calculated during pre-processing (260 and 878,**
 480 **respectively). The secondary mode was removed throughout the processing highlighting the effectiveness**
 481 **of common ion filtering.**

482

483 The main processing stage used a sample-to-blank ratio of 5 for blank subtraction with the
 484 mass-based approach. Allowable carbon ratios were set to $0.3 \leq H/C \leq 2.5$, $0 \leq O/C \leq 2$,
 485 $N/C \leq 1.3$, $S/C \leq 0.8$ and natural abundances were used for the carbon and sulphur isotopic
 486 ratios. The results of the main processing stage are shown in Figure 8 for both ionisation modes.
 487 A clear reduction in peaks is visible which is largely due to noise removal and blank
 488 subtraction, along with the additional filters, bringing the total assignment counts from 87,217
 489 and 238,006 to 720 and 6,491 for the APPI positive and ESI negative modes, respectively. The
 490 assignment count at this stage is equivalent to the peak count as duplicates were removed. Any
 491 peaks without assignments were removed from the spectra prior to analysis. The total
 492 assignment counts at the various stages of the data processing procedure, discussed further
 493 below, are summarised in Table 2. These values are specific to the current sample analysis and
 494 may not be representative for different analyses especially when adjusting the number of
 495 allowable formula assignments.

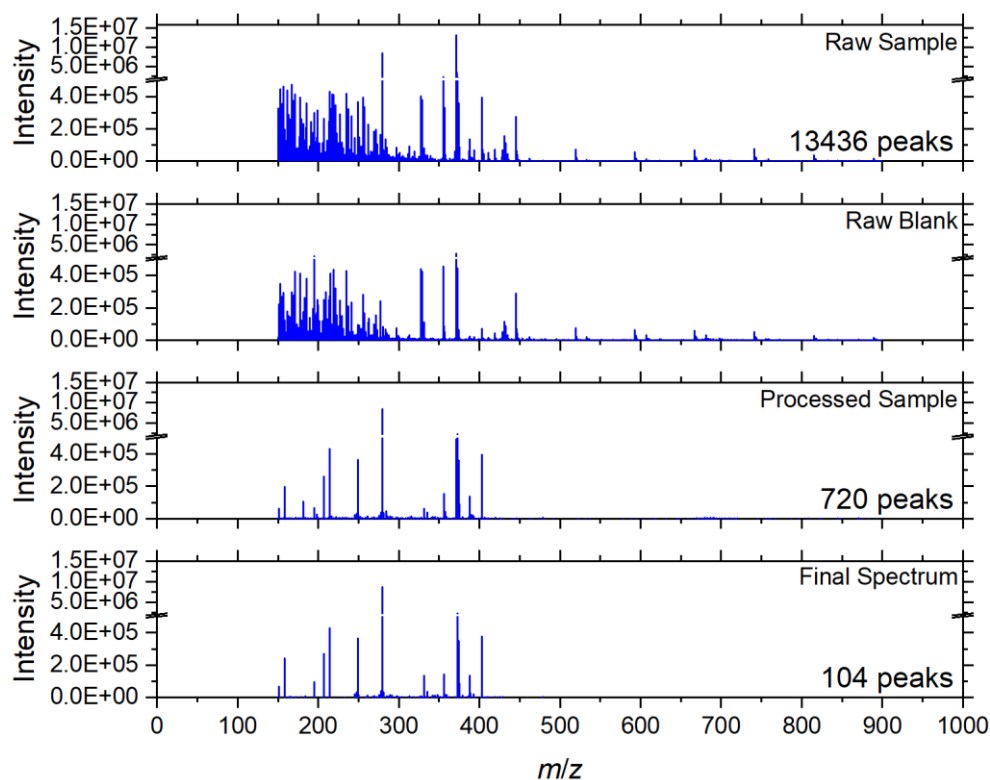
496 **Table 2** Number of assignments at various stages of processing, following the breakdown in Figure 1, given
 497 in absolute (Abs.) and relative (Rel.) terms. The assignment counts prior to common ion selection were
 498 based on a single repeat (same as in Figures 7 and 8). Common ion selection was based on four repeats for
 499 each ionisation mode. Duplicate assignments for a given peak were removed. These results are for the
 500 sample application only and may not be representative for different samples. The final number of
 501 assignments remaining correspond to 0.8% and 4.8% of the initial number of peaks for APPI+ and ESI-
 502 respectively.

Completed processing		APPI+		ESI-	
		Abs.	Rel.	Abs.	Rel.
Pre-treatment		87,217	100.0%	238,006	100.0%
Main processing	<i>Noise removal</i>	33,247	38.1%	123,757	53.3%
	<i>Blank subtraction</i>	17,468	20.0%	100,343	43.2%
	<i>Filter mass error range</i>	9,086	10.4%	78,999	34.0%
	<i>Filter carbon ratios</i>	5,311	6.1%	46,217	19.9%
	<i>DBE and nitrogen rule</i>	2,978	3.4%	23,812	10.3%
	<i>Isotopic filtering</i>	1,006	1.2%	7,681	3.3%
	<i>Duplicate removal</i> ^a	720	0.8%	6,491	2.8%
Common ion selection of four repeats		105	0.1%	2,603	1.1%
Shoulder ion removal		104	0.1%	2,598	1.1%

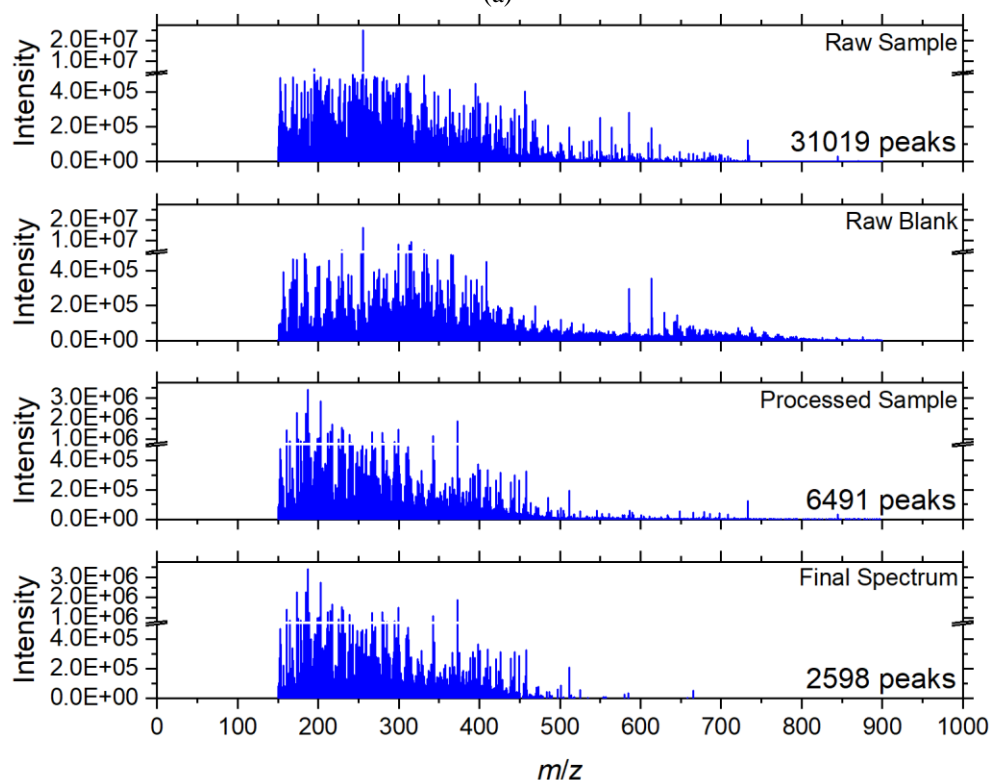
^a Removal of duplicate assignments is an optional processing step.

503
504

505 The four repeats of each ionisation mode were processed under the same conditions and
 506 subsequently processed to check for common ions. The common ion processing specified that
 507 a peak must appear in all four repeats to remain. This stage left 105 and 2,603 peaks for the
 508 APPI positive and ESI negative modes, respectively, before being checked for shoulder ions.
 509 Peaks were considered to be shoulders if their intensity was less than 1% of a high intensity
 510 peak (>1,000,000 a.u.) within 0.01 *m/z*. Only one and five peaks were considered shoulders for
 511 the APPI positive and ESI negative modes, respectively, largely due to the shoulder peaks
 512 being removed in earlier processing steps. The final spectrum for each ionisation mode is
 513 presented as the final panel in Figure 8. A more detailed description of shoulder ion removal is
 514 included in the Supplementary Material (Figure S2). The processing of the two sample sets
 515 concludes with 104 and 2,598 non-duplicate assignments remaining for the APPI positive and
 516 ESI negative modes, respectively.



(a)



(b)

517 **Figure 8** Mass spectra of the raw sample (first panel), blank (second panel), and processed sample (third
 518 panel) of a single replicate for the (a) APPI+ and (b) ESI- ionisation modes. The fourth panel is the final
 519 spectrum across four replicates after common ion retention and shoulder ion removal. Total peak numbers
 520 decreased from 13,436 to 104 and 31,019 to 2,598 for APPI+ and ESI-, respectively, with duplicate
 521 assignments removed. Note the raw sample peaks had multiple assignments (see Table 2).

522 **4 Conclusions**

523 In order to obtain realistic formula assignment from UHRMS data we introduced a processing
524 scheme based on knowledge of the samples being analysed and the instrument itself. The
525 developed scheme can be implemented for ESI and APPI ionisation techniques in both positive
526 and negative modes.

527 The procedure implements several common steps including *a priori* element selection and
528 exclusion filtering. Exclusion filtering, typically based on prior knowledge, includes known
529 instrument errors, general chemical principles (*e.g.* the nitrogen rule), and assumptions on
530 realistic element ratios. Pre-treatment considerations are also included to determine the noise
531 level and mass shift across each averaged spectrum.

532 Two methods of blank subtraction are available based on either processing the sample and
533 blank independently before comparison or performing the subtraction during the main
534 processing stage. While the former approach is more common, the latter approach typically
535 results in fewer false positives. The final spectrum also undergoes common ion selection to
536 exclude chemical noise peaks, when at least three replicates are acquired, and shoulder ion
537 removal for high intensity peaks.

538 A sample application of the processing scheme was presented, using both ionisation
539 techniques, to highlight the effectiveness of each stage in terms of peak removal. The final peak
540 counts were 0.8% and 8.4% for APPI positive and ESI negative ionisation, respectively,
541 relative to the number of peaks in the initial raw spectra with duplicate peak assignments
542 removed.

543 Potential future improvement of the overall processing scheme may involve the inclusion of
544 *m/z*-dependent mass shifts and additional elemental assignments such as phosphorus. The
545 current iteration, however, has already shown to be capable of processing complex atmospheric
546 compositions [12,22]. The approach allows for significant user-input enabling a wide range of
547 potential sample compositions and sampling methods (*e.g.* liquid extractive surface analysis of
548 flower petals [26]).

549 **Acknowledgement**

550 This work was supported by the European Research Council (ERC starting grant 279405) and
551 by the U.K. Natural Environment Research Council (NERC grant NE/H52449X/1). ATZ
552 thanks the Natural Sciences and Engineering Research Council of Canada, the Sir Winston
553 Churchill Society of Edmonton, and the Cambridge Trust for PhD funding. IK was supported
554 by a M. Curie Intra-European fellowship (project no. 254319).

555 **Supplementary material**

556 The codes associated with the processing scheme – named “Direct Infusion Mass Spectrometry
557 Processing (DIMaSP)” – are available for download at <https://doi.org/10.17863/CAM.9495>,
558 together with a basic manual. The current versions were tested for use with Wolfram
559 Mathematica 11.1.

560 **References**

- 561 [1] T. Kind, O. Fiehn, Seven golden rules for heuristic filtering of molecular formulas obtained by accurate
562 mass spectrometry, *BMC Bioinformatics*. 8 (2007) 105. doi:10.1186/1471-2105-8-105.
- 563 [2] I. Kourtchev, I.P. O'Connor, C. Giorio, S.J. Fuller, K. Kristensen, W. Maenhaut, J.C. Wenger, J.R.
564 Sodeau, M. Glasius, M. Kalberer, Effects of anthropogenic emissions on the molecular composition of
565 urban organic aerosols: an ultrahigh resolution mass spectrometry study, *Atmos. Environ.* 89 (2014) 525–
566 532. doi:10.1016/j.atmosenv.2014.02.051.
- 567 [3] S.A. Nizkorodov, J. Laskin, A. Laskin, Molecular chemistry of organic aerosols through the application
568 of high resolution mass spectrometry, *Phys. Chem. Chem. Phys.* 13 (2011) 3612–3629.
569 doi:10.1039/c0cp02032j.
- 570 [4] B.P. Koch, T. Dittmar, M. Witt, G. Kattner, Fundamentals of molecular formula assignment to ultrahigh
571 resolution mass data of natural organic matter, *Anal. Chem.* 79 (2007) 1758–1763.
572 doi:10.1021/ac061949s.
- 573 [5] G. Danger, F.-R. Orthous-Daunay, P. de Marcellus, P. Modica, V. Vuitton, F. Duvernay, L. Flandinet, L.
574 Le Sergeant d'Hendecourt, R. Thissen, T. Chiavassa, Characterization of laboratory analogs of
575 interstellar/cometary organic residues using very high resolution mass spectrometry, *Geochim.*
576 *Cosmochim. Acta.* 118 (2013) 184–201. doi:10.1016/j.gca.2013.05.015.
- 577 [6] C.C.L. Wong, D. Cociorva, J.D. Venable, T. Xu, J.R. Yates, Comparison of different signal thresholds on
578 data dependent sampling in Orbitrap and LTQ mass spectrometry for the identification of peptides and
579 proteins in complex mixtures, *J. Am. Soc. Mass Spectrom.* 20 (2009) 1405–1414.
580 doi:10.1016/j.jasms.2009.04.007.
- 581 [7] I. Kourtchev, S.J. Fuller, C. Giorio, R.M. Healy, E. Wilson, I. O'Connor, J.C. Wenger, M. McLeod, J.
582 Aalto, T.M. Ruuskanen, W. Maenhaut, R. Jones, D.S. Venables, J.R. Sodeau, M. Kulmala, M. Kalberer,
583 Molecular composition of biogenic secondary organic aerosols using ultrahigh-resolution mass
584 spectrometry: comparing laboratory and field studies, *Atmos. Chem. Phys.* 14 (2014) 2155–2167.
585 doi:10.5194/acp-14-2155-2014.
- 586 [8] B.P. Koch, M. Witt, R. Engbrodt, T. Dittmar, G. Kattner, Molecular formulae of marine and terrigenous
587 dissolved organic matter detected by electrospray ionization Fourier transform ion cyclotron resonance
588 mass spectrometry, *Geochim. Cosmochim. Acta.* 69 (2005) 3299–3308. doi:10.1016/j.gca.2005.02.027.
- 589 [9] J.W. Bright, E.C.M. Chen, Mass spectral interpretation using the “rule of ‘13,’” *J. Chem. Educ.* 60 (1983)
590 557. doi:10.1021/ed060p557.
- 591 [10] A.S. Wozniak, J.E. Bauer, R.L. Sleighter, R.M. Dickhut, P.G. Hatcher, Technical Note: Molecular
592 characterization of aerosol-derived water soluble organic carbon using ultrahigh resolution electrospray
593 ionization Fourier transform ion cyclotron resonance mass spectrometry, *Atmos. Chem. Phys.* 8 (2008)
594 5099–5111. doi:10.5194/acp-8-5099-2008.
- 595 [11] T. Ohno, P.E. Ohno, Influence of heteroatom pre-selection on the molecular formula assignment of soil
596 organic matter components determined by ultrahigh resolution mass spectrometry, *Anal. Bioanal. Chem.*
597 405 (2013) 3299–3306. doi:10.1007/s00216-013-6734-3.
- 598 [12] I. Kourtchev, S. Fuller, J. Aalto, T.M. Ruuskanen, M.W. McLeod, W. Maenhaut, R. Jones, M. Kulmala,
599 M. Kalberer, Molecular composition of boreal forest aerosol from Hyytiälä, Finland, using ultrahigh
600 resolution mass spectrometry, *Environ. Sci. Technol.* 47 (2013) 4069–4079. doi:10.1021/es3051636.
- 601 [13] Y. Lu, X. Li, R. Mesfioui, J.E. Bauer, R.M. Chambers, E.A. Canuel, P.G. Hatcher, Use of ESI-FTICR-
602 MS to characterize dissolved organic matter in headwater streams draining forest-dominated and pasture-
603 dominated watersheds, *PLoS One.* 10 (2015) 1–21. doi:10.1371/journal.pone.0145639.
- 604 [14] E.B. Kujawinski, M.D. Behn, Automated analysis of electrospray ionization fourier transform ion
605 cyclotron resonance mass spectra of natural organic matter, *Anal. Chem.* 78 (2006) 4363–4373.
606 doi:10.1021/ac0600306.
- 607 [15] T. Ohno, R.L. Sleighter, P.G. Hatcher, Comparative study of organic matter chemical characterization
608 using negative and positive mode electrospray ionization ultrahigh-resolution mass spectrometry, *Anal.*
609 *Bioanal. Chem.* 408 (2016) 2497–2504. doi:10.1007/s00216-016-9346-x.
- 610 [16] G. Glauser, N. Veyrat, B. Rochat, J.-L. Wolfender, T.C.J. Turlings, Ultra-high pressure liquid
611 chromatography-mass spectrometry for plant metabolomics: a systematic comparison of high-resolution
612 quadrupole-time-of-flight and single stage Orbitrap mass spectrometers, *J. Chromatogr. A.* 1292 (2013)
613 151–159. doi:10.1016/j.chroma.2012.12.009.
- 614 [17] C.A. Hughey, C.L. Hendrickson, R.P. Rodgers, A.G. Marshall, K. Qian, Kendrick mass defect spectrum:
615 a compact visual analysis for ultrahigh-resolution broadband mass spectra, *Anal. Chem.* 73 (2001) 4676–
616 4681. doi:10.1021/ac010560w.
- 617 [18] A.L. Putman, J.H. Offenberg, R. Fisseha, S. Kundu, T.A. Rahn, L.R. Mazzoleni, Ultrahigh-resolution FT-
618 ICR mass spectrometry characterization of α -pinene ozonolysis SOA, *Atmos. Environ.* 46 (2012) 164–

- 619 172. doi:10.1016/j.atmosenv.2011.10.003.
- 620 [19] J. V Olsen, L.M.F. de Godoy, G. Li, B. Macek, P. Mortensen, R. Pesch, A. Makarov, O. Lange, S.
621 Horning, M. Mann, Parts per million mass accuracy on an Orbitrap mass spectrometer via lock mass
622 injection into a C-trap, *Mol. Cell. Proteomics*. 4 (2005) 2010–2021. doi:10.1074/mcp.T500030-MCP200.
- 623 [20] T. Reemtsma, Determination of molecular formulas of natural organic matter molecules by (ultra-) high-
624 resolution mass spectrometry: status and needs, *J. Chromatogr. A*. 1216 (2009) 3687–3701.
625 doi:10.1016/j.chroma.2009.02.033.
- 626 [21] P. Herzprung, N. Hertkorn, W. von Tümpling, M. Harir, K. Friese, P. Schmitt-Kopplin, Molecular
627 formula assignment for dissolved organic matter (DOM) using high-field FT-ICR-MS: chemical
628 perspective and validation of sulphur-rich organic components (CHOS) in pit lake samples, *Anal. Bioanal.*
629 *Chem.* 408 (2016) 2461–2469. doi:10.1007/s00216-016-9341-2.
- 630 [22] I. Kourtchev, R.H.M. Godoi, S. Connors, J.G. Levine, A.T. Archibald, A.F.L. Godoi, S.L. Paralovo,
631 C.G.G. Barbosa, R.A.F. Souza, A.O. Manzi, R. Seco, S. Sjostedt, J.-H. Park, A. Guenther, S. Kim, J.
632 Smith, S.T. Martin, M. Kalberer, Molecular composition of organic aerosols in central Amazonia: an
633 ultra-high-resolution mass spectrometry study, *Atmos. Chem. Phys.* 16 (2016) 11899–11913.
634 doi:10.5194/acp-16-11899-2016.
- 635 [23] E. V Kunenkov, A.S. Kononikhin, I. V Perminova, N. Hertkorn, A. Gaspar, P. Schmitt-Kopplin, I.A.
636 Popov, A. V Garmash, E.N. Nikolaev, Total mass difference statistics algorithm: a new approach to
637 identification of high-mass building blocks in electrospray ionization Fourier transform ion cyclotron
638 mass spectrometry data of natural organic matter, *Anal. Chem.* 81 (2009) 10106–10115.
639 doi:10.1021/ac901476u.
- 640 [24] T.B. Nguyen, A.P. Bateman, D.L. Bones, S.A. Nizkorodov, J. Laskin, A. Laskin, High-resolution mass
641 spectrometry analysis of secondary organic aerosol generated by ozonolysis of isoprene, *Atmos. Environ.*
642 44 (2010) 1032–1042. doi:10.1016/j.atmosenv.2009.12.019.
- 643 [25] E. de Hoffmann, V. Stroobant, *Mass spectrometry: principles and applications*, 3rd ed., John Wiley &
644 Sons, Ltd., Chichester, England, 2007.
- 645 [26] C. Giorio, E. Moyroud, B.J. Glover, P.C. Skelton, M. Kalberer, Direct surface analysis coupled to high-
646 resolution mass spectrometry reveals heterogeneous composition of the cuticle of *Hibiscus trionum* petals,
647 *Anal. Chem.* 87 (2015) 9900–9907. doi:10.1021/acs.analchem.5b02498.
- 648 [27] P. Schmitt-Kopplin, A. Gelencsér, E. Dabek-Zlotorzynska, G. Kiss, N. Hertkorn, M. Harir, Y. Hong, I.
649 Gebefügi, Analysis of the unresolved organic fraction in atmospheric aerosols with ultrahigh-resolution
650 mass spectrometry and nuclear magnetic resonance spectroscopy: organosulfates as photochemical smog
651 constituents, *Anal. Chem.* 82 (2010) 8017–8026. doi:10.1021/ac101444r.
- 652 [28] L.R. Mazzoleni, B.M. Ehrmann, X. Shen, A.G. Marshall, J.L. Collett, Water-soluble atmospheric organic
653 matter in fog: Exact masses and chemical formula identification by ultrahigh-resolution fourier transform
654 ion cyclotron resonance mass spectrometry, *Environ. Sci. Technol.* 44 (2010) 3690–3697.
655 doi:10.1021/es903409k.
- 656 [29] C. Ranninger, L.E. Schmidt, M. Rurik, A. Limonciel, P. Jennings, O. Kohlbacher, C.G. Huber, Improving
657 global feature detectabilities through scan range splitting for untargeted metabolomics by high-
658 performance liquid chromatography-Orbitrap mass spectrometry, *Anal. Chim. Acta.* 930 (2016) 13–22.
659 doi:10.1016/j.aca.2016.05.017.
- 660 [30] R.L. Sleighter, H. Chen, A.S. Wozniak, A.S. Willoughby, P. Caricasole, P.G. Hatcher, Establishing a
661 measure of reproducibility of ultrahigh-resolution mass spectra for complex mixtures of natural organic
662 matter, *Anal. Chem.* 84 (2012) 9184–9191. doi:10.1021/ac3018026.
- 663 [31] A. Makarov, E. Denisov, O. Lange, S. Horning, Dynamic range of mass accuracy in LTQ Orbitrap hybrid
664 mass spectrometer, *J. Am. Soc. Mass Spectrom.* 17 (2006) 977–982. doi:10.1016/j.jasms.2006.03.006.
- 665 [32] S. Kim, R.P. Rodgers, A.G. Marshall, Truly “exact” mass: elemental composition can be determined
666 uniquely from molecular mass measurement at ~0.1 mDa accuracy for molecules up to ~500 Da, *Int. J.*
667 *Mass Spectrom.* 251 (2006) 260–265. doi:10.1016/j.ijms.2006.02.001.
- 668 [33] A. Tapparo, C. Giorio, L. Soldà, S. Bogianni, D. Marton, M. Marzaro, V. Girolami, UHPLC-DAD method
669 for the determination of neonicotinoid insecticides in single bees and its relevance in honeybee colony
670 loss investigations, *Anal. Bioanal. Chem.* 405 (2013) 1007–1014. doi:10.1007/s00216-012-6338-3.
- 671 [34] M.C.K. Soule, K. Longnecker, S.J. Giovannoni, E.B. Kujawinski, Impact of instrument and experiment
672 parameters on reproducibility of ultrahigh resolution ESI FT-ICR mass spectra of natural organic matter,
673 *Org. Geochem.* 41 (2010) 725–733. doi:10.1016/j.orggeochem.2010.05.017.
- 674 [35] J.A. Hawkes, T. Dittmar, C. Patriarca, L. Tranvik, J. Bergquist, Evaluation of the Orbitrap mass
675 spectrometer for the molecular fingerprinting analysis of natural dissolved organic matter, *Anal. Chem.*
676 88 (2016) 7698–7704. doi:10.1021/acs.analchem.6b01624.
- 677 [36] R.E. O’Brien, A. Laskin, J. Laskin, S. Liu, R. Weber, L.M. Russell, A.H. Goldstein, Molecular
678 characterization of organic aerosol using nanospray desorption/electrospray ionization mass

679 spectrometry: CalNex 2010 field study, *Atmos. Environ.* 68 (2013) 265–272.
680 doi:10.1016/j.atmosenv.2012.11.056.

681 [37] D. Freedman, P. Diaconis, On the histogram as a density estimator: L2 theory, *Zeitschrift Für*
682 *Wahrscheinlichkeitstheorie Und Verwandte Gebiete.* 57 (1981) 453–476. doi:10.1007/BF01025868.

683 [38] K.O. Zhurov, A.N. Kozhinov, L. Fornelli, Y.O. Tsybin, Distinguishing analyte from noise components in
684 mass spectra of complex samples: Where to cut the noise?, *Anal. Chem.* 86 (2014) 3308–3316.
685 doi:10.1021/ac403278t.

686 [39] S.J. Fuller, Y. Zhao, S.S. Cliff, A.S. Wexler, M. Kalberer, Direct surface analysis of time-resolved aerosol
687 impactor samples with ultrahigh-resolution mass spectrometry, *Anal. Chem.* 84 (2012) 9858–9864.
688 doi:10.1021/ac3020615.

689 [40] R.L. Sleighter, G.A. McKee, Z. Liu, P.G. Hatcher, Naturally present fatty acids as internal calibrants for
690 Fourier transform mass spectra of dissolved organic matter, *Limnol. Oceanogr. Methods.* 6 (2008) 246–
691 253. doi:10.4319/lom.2008.6.246.

692 [41] A.N. Kozhinov, K.O. Zhurov, Y.O. Tsybin, Iterative method for mass spectra recalibration via empirical
693 estimation of the mass calibration function for fourier transform mass spectrometry-based petroleomics,
694 *Anal. Chem.* 85 (2013) 6437–6445. doi:10.1021/ac400972y.

695 [42] F.E. Grubbs, Procedures for Detecting Outlying Observations in Samples, *Technometrics.* 11 (1969) 1–
696 21. doi:10.1080/00401706.1969.10490657.

697 [43] X.K. Wang, S. Rossignol, Y. Ma, L. Yao, M.Y. Wang, J.M. Chen, C. George, L. Wang, Molecular
698 characterization of atmospheric particulate organosulfates in three megacities at the middle and lower
699 reaches of the Yangtze River, *Atmos. Chem. Phys.* 16 (2016) 2285–2298. doi:10.5194/acp-16-2285-2016.

700 [44] A.G. Rincón, A.I. Calvo, M. Dietzel, M. Kalberer, Seasonal differences of urban organic aerosol
701 composition - an ultra-high resolution mass spectrometry study, *Environ. Chem.* 9 (2012) 298–319.
702 doi:10.1071/EN12016.

703 [45] A.P. Bateman, S.A. Nizkorodov, J. Laskin, A. Laskin, Time-resolved molecular characterization of
704 limonene/ozone aerosol using high-resolution electrospray ionization mass spectrometry, *Phys. Chem.*
705 *Chem. Phys.* 11 (2009) 7931–7942. doi:10.1039/b905288g.

706 [46] R.L. Sleighter, P.G. Hatcher, Molecular characterization of dissolved organic matter (DOM) along a river
707 to ocean transect of the lower Chesapeake Bay by ultrahigh resolution electrospray ionization Fourier
708 transform ion cyclotron resonance mass spectrometry, *Mar. Chem.* 110 (2008) 140–152.
709 doi:10.1016/j.marchem.2008.04.008.

710 [47] B.P. Koch, T. Dittmar, From mass to structure: an aromaticity index for high-resolution mass data of
711 natural organic matter, *Rapid Commun. Mass Spectrom.* 20 (2006) 926–932. doi:10.1002/rcm.2386.

712 [48] K.E. Altieri, B.J. Turpin, S.P. Seitzinger, Oligomers, organosulfates, and nitrooxy organosulfates in
713 rainwater identified by ultra-high resolution electrospray ionization FT-ICR mass spectrometry, *Atmos.*
714 *Chem. Phys.* 9 (2009) 2533–2542. doi:10.5194/acp-9-2533-2009.

715 [49] A. Stubbins, R.G.M. Spencer, H. Chen, P.G. Hatcher, K. Mopper, P.J. Hernes, V.L. Mwamba, A.M.
716 Mangangu, J.N. Wabakanghanzi, J. Six, Illuminated darkness: molecular signatures of Congo River
717 dissolved organic matter and its photochemical alteration as revealed by ultrahigh precision mass
718 spectrometry, *Limnol. Oceanogr.* 55 (2010) 1467–1477. doi:10.4319/lo.2010.55.4.1467.

719 [50] D.R. Fooshee, P.K. Aiona, A. Laskin, J. Laskin, S.A. Nizkorodov, P.F. Baldi, Atmospheric oxidation of
720 squalene: molecular study using COBRA modeling and high-resolution mass spectrometry, *Environ. Sci.*
721 *Technol.* 49 (2015) 13304–13313. doi:10.1021/acs.est.5b03552.

722 [51] J.R. de Laeter, J.K. Böhlke, P. De Bièvre, H. Hidaka, H.S. Peiser, K.J.R. Rosman, P.D.P. Taylor, Atomic
723 weights of the elements: review 2000 (IUPAC technical report), *Pure Appl. Chem.* 75 (2003) 785.
724 doi:10.1351/pac200375060683.

725 [52] Y. Xu, J.-F. Heilier, G. Madalinski, E. Genin, E. Ezan, J.-C. Tabet, C. Junot, Evaluation of accurate mass
726 and relative isotopic abundance measurements in the LTQ-Orbitrap mass spectrometer for further
727 metabolomics database building, *Anal. Chem.* 82 (2010) 5490–5501. doi:10.1021/ac100271j.

728 [53] P. Herzsprung, N. Hertkorn, W. von Tümpling, M. Harir, K. Friese, P. Schmitt-Kopplin, Understanding
729 molecular formula assignment of Fourier transform ion cyclotron resonance mass spectrometry data of
730 natural organic matter from a chemical point of view, *Anal. Bioanal. Chem.* 406 (2014) 7977–7987.
731 doi:10.1007/s00216-014-8249-y.

732 [54] T. Pluskal, S. Castillo, A. Villar-Briones, M. Oresic, MZmine 2: modular framework for processing,
733 visualizing, and analyzing mass spectrometry-based molecular profile data, *BMC Bioinformatics.* 11
734 (2010) 395. doi:10.1186/1471-2105-11-395.

735 [55] C. Giorio, A. Tapparo, M.L. Scapellato, M. Carrieri, P. Apostoli, G.B. Bartolucci, Field comparison of a
736 personal cascade impactor sampler, an optical particle counter and CEN-EU standard methods for PM10,
737 PM2.5 and PM1 measurement in urban environment, *J. Aerosol Sci.* 65 (2013) 111–120.
738 doi:10.1016/j.jaerosci.2013.07.013.