

Cite this: DOI: 10.1039/xxxxxxxxxx

# On the benefits of using multivariate analysis in mass spectrometric studies of combustion-generated aerosols

D. Duca<sup>a</sup>, C. Irimiea<sup>b</sup>, A. Faccinetto<sup>c</sup>, J. A. Noble<sup>a,†</sup>, M. Vojkovic<sup>a</sup>, Y. Carpentier<sup>a</sup>, I. K. Ortega<sup>b</sup>, C. Pirim<sup>a</sup> and C. Focsa<sup>a</sup>

Received Date

Accepted Date

DOI: 10.1039/xxxxxxxxxx

www.rsc.org/journalname

Detailed molecular-level analysis of combustion emissions may be challenging even with high-resolution mass spectrometry. The intricate chemistry of the carbonaceous particles surface layer (which drives their reactivity, environmental and health impacts) results in complex mass spectra. Building on a recently proposed comprehensive methodology (encompassing all stages from sampling to data reduction), we propose herein a comparative analysis of soot particles produced by three different sources: a miniCAST standard generator, a laboratory diffusion flame and a single cylinder internal combustion engine. The surface composition is probed by either laser or secondary ion mass spectrometry. Principal component analysis and hierarchical clustering analysis proved their efficiency in both identifying general trends and evidencing subtle differences that otherwise would remain unnoticed in the plethora of data generated during mass spectrometric analyses. Chemical information extracted from these multivariate statistical procedures contributes to a better understanding of fundamental combustion processes and also opens to practical applications such as the tracing of engine emissions.

## 1 Introduction

Multivariate analysis (MVA) methods are powerful tools to unravel trends in complex databases. They have been successfully applied in the past, for instance, to identify drug metabolites in biological fluids<sup>1</sup>, to evaluate profiles of volatile compounds present in mainstream tobacco smoke<sup>2</sup>, or else, to assess surface water quality<sup>3</sup>. Among the MVA methods commonly used<sup>4</sup> are the principal component analysis (PCA) and the hierarchical clustering analysis (HCA). The former is used to reveal hidden patterns in databases, by emphasising the variance between samples and thus highlighting their differences and similarities<sup>5</sup>, whereas the latter searches for patterns in a database by grouping the observables into distinct clusters. Their capability at distinguishing various complex samples, as exemplified for a while now in the field of biology, has recently led to their consideration for unravelling the chemical composition of multifaceted samples of envi-

ronmental interest.

Atmospheric aerosols are airborne particles consisting of an intricate mixture of chemical constituents whose nature varies greatly depending upon their emission source and evolution within the atmosphere. Carbonaceous particles account for a significant fraction of atmospheric particulate matter in urban areas (typically 30-50% by mass<sup>6-8</sup>). They are mainly formed of soot, i.e. particles generated by the incomplete combustion of hydrocarbon-based fuels or biomass. Accordingly, soot particles hold a multitude of chemical compounds derived from various sources (remnant of fuels, combustion and/or post-oxidation products, etc.) that may have been further transformed (aged) by the time they are analysed due to their continuous interaction with environmental elements (solar rays, water molecules, pollutants, etc.). Soot particles are therefore considered a complex mixture that often needs a concerted analytical scheme to be fully resolved.

Mass spectrometry (MS) based techniques have significantly contributed to better understanding soot chemistry over the years. They are generally robust techniques that do not require extensive sample preparation, and are hence preferred for the analysis of such complex samples. Furthermore, the amount of particulate matter required to perform MS analysis is relatively small. MS based techniques mostly differ by the way the ions

<sup>a</sup> Univ. Lille, CNRS, UMR 8523 – PhLAM – Laboratoire de Physique des Lasers Atomes et Molécules, F-59000 Lille, France. E-mail: dimitru.duca@univ-lille.fr

<sup>b</sup> ONERA – The French Aerospace Laboratory, F-91123 Palaiseau, France.

<sup>c</sup> Univ. Lille, CNRS, UMR 8522 – PC2A – Laboratoire de Physico-Chimie des Processus de Combustion de l'Atmosphère, F-59000 Lille, France.

<sup>†</sup> Current address: CNRS, Aix Marseille Université, PIIM, UMR 7345, 13397 Marseille cedex, France.

transferred to the mass spectrometer are created (e.g. soot particle aerosol mass spectrometry (SP-AMS)<sup>9</sup>, two-step laser mass spectrometry (L2MS)<sup>10</sup>, time-of-flight secondary ion mass spectrometry (ToF-SIMS)<sup>11,12</sup>), which often condition their specificity to provide information on either bulk or surface chemical composition. Ultra high resolution mass analyzers as Orbitrap, Fourier transform ion cyclotron resonance (FT-ICR) and high resolution quadrupole time of flight MS can reach a resolving power higher than 90 000<sup>13,14</sup>. These techniques were developed mainly for proteomics and pharmaceutical analyses, but lately their application has been extended to many other fields among which they start being used and adapted to atmospheric aerosols<sup>15,16</sup>. However, ultra high resolution mass spectrometry is still very rarely applied to the analysis of combustion products, with only a few examples to date<sup>17</sup>. Ultra high resolution mass analyzers are powerful analytical tools, however they still need validation of the sampling protocols. For instance, the sample transfer into the instrument is based on nanospray desorption electrospray using a polar solvent for Orbitrap, followed more recently by laser desorption for FT-ICR and atmospheric pressure chemical ionization (APCI) for APCI-Orbitrap<sup>13,16,17</sup>. Let us also emphasize that in directed energy (laser, ion beam) desorption methods, beside the analyzer performances, the condensed-gas phase transfer itself plays a critical role in the maximum achievable mass resolution and on the total number of detected signals, through, e.g., the sample/substrate roughness or conductive properties. We therefore stress the need for a thorough evaluation (and optimization) of the entire analysis chain, from sample collection/deposition on suitable substrates, to sample transfer/ionization into gas phase, ions mass separation and detection, and finally powerful data treatment and interpretation<sup>18,19</sup>.

Mass spectra of soot particles can be very complex, featuring hundreds and even thousands mass peaks, which quickly renders the interpretation of mass spectra difficult and therefore limits the potentiality of MS to resolve complex mixtures. Accordingly, resolving sample complexity in MS databases is currently tackled using two main approaches. The first is based on the identification of marker species, i.e. compounds that are directly linked to a source/process and that can thus be considered as their fingerprints, while the second approach relies on statistical methods. In particular, the use of MVA methods in conjunction with MS is a creative combination to exploit all of the information given by a multitude of peaks within a great variety of sample sets. Both approaches are widely used in analysis of mass spectra obtained with aerosol mass spectrometers (AMS)<sup>20–22</sup>, proton transfer reaction mass spectrometers (PTR-MS)<sup>23,24</sup>, and laser-based MS techniques<sup>19,25,26</sup>. Discrimination using marker species was applied to samples of various sources, proving its effectiveness when comparing soot emitted from wood combustion<sup>20,27</sup>, on-road vehicles<sup>25</sup>, aircrafts<sup>22–24,28,29</sup>, ships<sup>30</sup> or other ambient aerosols<sup>21</sup>. However, since some marker species may not remain stable over the aerosols' life span, especially upon atmospheric ageing<sup>6</sup>, this method may misdirect with regards to the origin of samples a priori unknown. To circumvent this limitation, MVA approaches are chosen, as they can discriminate samples regardless of their provenance or evolution. Therefore, MVA can uncover trends and

features even in samples of unknown/mixed origins<sup>28,31</sup>, which is particularly interesting when analysing natural aerosols.

In constant interaction with their surroundings, aerosols surfaces drive their overall reactivity, and therefore, set their evolution path within the atmosphere (sedimentation, formation of secondary organic aerosols, nucleation, etc.). It is hence imperative to uncover their complex surface composition in order to assess their impact on both human health and the environment<sup>32,33</sup>. For example, some polycyclic aromatic hydrocarbons (PAHs), often found adsorbed on the surface of soot particles, are known to be toxic and to have mutagenic effects<sup>34,35</sup>. In addition, the chemical composition of aerosol surfaces determines their hygroscopicity<sup>36</sup> and therefore their ability to act as condensation nuclei, potentially influencing climate forcing, cloud cover and precipitations.

Our group has been addressing this issue of untangling surface chemical compositions of field-collected or laboratory-generated combustion aerosols for over a decade<sup>10,18,19,26,29,30,37–40</sup>. We recently described an original and comprehensive experimental methodology<sup>18</sup> that we later implemented in combining statistical-based approaches with compound classification techniques<sup>19</sup>. This latter systematic study by Irimiea and coworkers<sup>19</sup> was undertaken to characterise over 100 samples collected from different flames. In this work we developed a comprehensive protocol that allowed significant progress towards the fundamental understanding of soot nucleation and growth. Laboratory flames or standard soot generators are often used to produce soot particles with similar physico-chemical properties to the ones produced by “real world” combustion sources<sup>41</sup>. Laboratory soot particles offer the advantages of a reproducible, easy-access and low-cost production, which is of great importance when testing the robustness of a protocol. Therefore, this necessary step is of paramount importance for further refinements in field-collected combustion-generated particle analyses.

## 2 Experimental

In this section, the choice of the combustion conditions, the sampling approach and the experimental techniques used to characterise the samples are detailed. In particular, L2MS and SIMS are used in parallel to obtain information on the chemical composition of combustion generated aerosols.

### 2.1 Soot samples

Soot samples are generated in different combustion conditions (fuel, burner and sampling method) in order to test the ability of our data treatment protocols to reveal differences and similarities between samples. The sampling procedure, including the substrate choice and its preparation, is optimised according to our previous experience<sup>18</sup>. In particular, the sample-substrate reactivity can lead to the formation of a large number of byproducts that clutter the mass spectrum and make the identification of individual compounds much more difficult. A short description of all analysed samples (summarised in Table 1) is given below. The following soot samples have been used:

- Soot produced by a miniCAST generator (5201c) from Jing

150 Ltd., which is currently proposed as a means of obtaining<sup>200</sup>  
151 “standard” soot easily comparable to other studies<sup>41–43</sup>. The<sup>201</sup>  
152 main difference between the miniCAST working points is the<sup>202</sup>  
153 oxidation flow (1.50 → 1.15 → 1.00 L min<sup>-1</sup>) resulting in<sup>203</sup>  
154 three different combustion conditions (C<sub>1</sub> → C<sub>2</sub> → C<sub>3</sub>)<sup>41–43</sup>.<sup>204</sup>  
155 The hereby generated particles are subsequently deposited<sup>205</sup>  
156 on quartz fibre filters. <sup>206</sup>

157 • Soot produced by laboratory turbulent diffusion flames sup-<sup>208</sup>  
158 plied with two different liquid fuels: diesel (D1-5) and<sup>209</sup>  
159 kerosene (K1-5). Soot particles are sampled from the flame<sup>210</sup>  
160 at different height above the burner (HAB) and deposited by<sup>211</sup>  
161 impaction on Si wafers. Sampling at various HAB is a means<sup>212</sup>  
162 of investigating soot particles of different maturity<sup>38</sup>. <sup>213</sup>

163 • Soot produced by a gasoline single cylinder internal combus-<sup>215</sup>  
164 tion engine (ICE). Operating conditions of this engine (e.g.<sup>216</sup>  
165 injection and ignition crank angle, applied load) could be<sup>217</sup>  
166 easily changed, thus allowing exhausts sampling at various<sup>218</sup>  
167 working regimes. The following operating points were used:<sup>219</sup>

168 – normal engine operation, i.e. engine optimised in<sup>220</sup>  
169 terms of high efficiency and low particle emissions,<sup>221</sup>  
170 with medium (GOM) and high (GOH) applied loads,<sup>222</sup>  
171 which simulate different driving regimes; <sup>223</sup>  
172 <sup>224</sup>

173 – malfunction simulation with a medium load applied;<sup>225</sup>  
174 low air/fuel ratio resulting in a high-sooting regime<sup>226</sup>  
175 (GEF) and an addition of oil to the combustion cham-<sup>227</sup>  
176 ber (GEO). <sup>228</sup>  
177 <sup>229</sup>

178 Soot particles are sampled using a cascade impactor<sup>230</sup>  
179 (NanoMOUDI) to enable for size selection during sampling,<sup>231</sup>  
180 and deposited on Al foils. We analysed the particles collected<sup>232</sup>  
181 on the last five stages, having diameter in the range 10-180<sup>233</sup>  
182 nm (Table 1). <sup>234</sup>  
183 <sup>235</sup>

184 Off-line analysis of soot particles requires a careful choice of  
185 the deposition substrate, not only to minimise the risk of contam-  
186 inating the samples, but also to ensure that a high mass resolution<sup>236</sup>  
187 can be achieved. In particular, among other factors, the mass res-<sup>237</sup>  
188 olution is directly linked to the surface roughness of the substrate,<sup>238</sup>  
189 and can be maximised by depositing the samples on ultra-flat sur-<sup>239</sup>  
190 faces such as Si or Ti wafers. Furthermore, the sample-substrate<sup>240</sup>  
191 reactivity can lead to the formation of reaction byproducts that<sup>241</sup>  
192 may heavily interfere with the assignment of sample-specific sig-<sup>242</sup>  
193 nals. Therefore, the careful characterization/choice of the deposi-<sup>243</sup>  
194 tion substrate is mandatory and the comprehensive identification<sup>244</sup>  
195 of its possible reactivity byproducts is necessary for a valid analyt-<sup>245</sup>  
196 ical protocol<sup>18,19</sup>. Regardless of its nature, the substrate should<sup>246</sup>  
197 undergo a series of preparation steps before it can be used to col-  
198 lect particulate matter. <sup>247</sup>  
199 <sup>248</sup>

## 196 2.2 Two-step laser mass spectrometry (L2MS) <sup>249</sup>

197 This laser-based MS technique has been extensively used by our<sup>250</sup>  
198 group to characterise the chemical composition of combustion<sup>251</sup>  
199 byproducts during the last decade<sup>10,18,26,29,30,37–39</sup>. The main<sup>252</sup>

advantages of L2MS are its high sensitivity and selectivity with re-  
gards to specific classes of compounds thanks to resonant ionisa-  
tion processes that can be tuned to reach for instance the sub-fmol  
limit for the detection of PAHs<sup>10,37</sup>. In addition, the controlled  
laser desorption process ensures a soft removal of molecules ad-  
sorbed on the particle surface (typically sub-monolayer regime),  
and thus avoids/limits either their fragmentation or the in-depth  
damaging of the underlying carbon matrix<sup>37</sup>. This qualifies L2MS  
as a surface-sensitive analysis technique, comparable in limit of  
detection ( $\sim 10^{-6}$  monolayers) with static-mode secondary ions  
mass spectrometry (SIMS, see below), but with much lower an-  
alyte fragmentation. However, our previous L2MS studies were  
limited by a mass resolution of  $m/\Delta m \sim 1000$ , significantly lower  
than the one achievable in SIMS (up to  $m/\Delta m \sim 10\,000$ , depend-  
ing on the deposition substrate<sup>18,19</sup>). In the current work, we  
take benefit of the recent implementation of a new mass spec-  
trometer (Fasmatech S&T) which combines ion cooling, Radio  
Frequency (RF) guiding and Time of Flight (ToF) analyser to  
reach a mass resolution of about  $m/\Delta m \sim 15000$ . In this new ex-  
perimental setup, the sample, placed under vacuum ( $10^{-8}$  mbar  
residual pressure), is irradiated at 30° angle of incidence by a fre-  
quency doubled Nd:YAG laser beam (Quantel Brilliant,  $\lambda = 532$   
nm, 4 ns pulse duration,  $\sim 50$  mJ cm<sup>-2</sup> fluence, 10 Hz repetition  
rate) focused to a 0.3 mm<sup>2</sup> spot on the surface. The desorbed  
compounds form a gas plume expanding in the vacuum normally  
to the sample surface, and are ionised by an orthogonal UV laser  
beam (Quantel Brilliant,  $\lambda_i = 266$  nm, 4 ns pulse duration, 10  
Hz repetition rate,  $\sim 0.3$  J cm<sup>-2</sup> fluence). At this ionisation wave-  
length, a high sensitivity is achieved for PAHs through a resonance  
enhanced multiphoton ionisation process 1+1 REMPI<sup>44–46</sup>. Care  
must be taken on the coupling of the desorption and ionisation  
steps in this laser-based MS technique<sup>47–49</sup>. Moreover, by chang-  
ing the ionisation wavelength, one can target different classes of  
compounds. The generated ions are then RF-guided to a He colli-  
sion cell for thermalisation and subsequently mass analysed in a  
time of flight mass spectrometer (ToF-MS).

## 2.3 Secondary Ion Mass Spectrometry (SIMS)

In addition, the samples are characterised by using a commercial  
IONTOF ToF-SIMS<sup>5</sup> secondary ion mass spectrometer with maxi-  
mum resolving power of  $m/\Delta m \sim 10\,000$ . In short, samples are  
placed in the analysis chamber with a residual pressure of  $\sim 10^{-7}$   
mbar. The surface of the sample is bombarded by a 25 keV Bi<sup>3+</sup>  
ion beam with a current of 0.3 pA in static mode. A small fraction  
of the ejected atoms/molecules are ionised (secondary ions) and  
can thus be analysed using a time-of-flight tube (V mode). Mass  
spectra are recorded in both positive and negative polarities, to  
obtain the maximum amount of information on the sample<sup>18,19</sup>.

## 3 Data Analysis Methodology and Exam- ples of Applications

The data presented below is analysed following an approach  
structured in three main points that include: mass defect analysis  
for identification of unknown compounds (Section 3.1), multi-  
variate analysis for the reduction of the number of dimensions of

**Table 1** Soot samples used to put in evidence the proposed methodology

Name	Fuel	Source	Substrate	Description	Analysing technique
C1	propane	miniCAST	Quartz fibre filters	1.5 l/min oxidation flow	L2MS +
C2				1.15 l/min oxidation flow	
C3				1.0 l/min oxidation flow	
D1	diesel	diffusion flame	Si wafer	HAB = 6mm	SIMS +/-
D2				HAB = 12mm	
D3				HAB = 14mm	
D4				HAB = 18mm	
D5				HAB = 24mm	
K1	kerosene	diffusion flame	Si wafer	HAB = 6mm	SIMS +/-
K2				HAB = 12mm	
K3				HAB = 14mm	
K4				HAB = 18mm	
K5				HAB = 24mm	
GOM1	gasoline	ICE, optimal conditions, medium load	Al foil	∅100 - 180nm	SIMS +/-
GOM2				∅56 - 100nm	
GOM3				∅32 - 56nm	
GOM4				∅18 - 32 nm	
GOM5				∅10 - 18 nm	
GOH1	gasoline	ICE, optimal conditions, high load	Al foil	∅100 - 180nm	SIMS +/-
GOH2				∅56 - 100nm	
GOH3				∅32 - 56nm	
GOH4				∅18 - 32 nm	
GEF1	gasoline	ICE, low Air/Fuel ratio	Al foil	∅100 - 180nm	SIMS +/-
GEF2				∅56 - 100nm	
GEF3				∅32 - 56nm	
GEF4				∅18 - 32 nm	
GEO1	gasoline	ICE, addition of oil	Al foil	∅100 - 180nm	SIMS +/-
GEO2				∅56 - 100nm	
GEO3				∅32 - 56nm	
GEO4				∅18 - 32 nm	

the dataset (Section 3.2) and eventually mass peak grouping for uncovering hidden trends and highlight correlations between different classes of compounds (Section 3.3). This section details the proposed data treatment protocol. Mass spectra of the previously described samples have been used to demonstrate its advantages, including its universal character (the ability to be used with mass spectra of various samples, obtained with different experimental techniques). Mass spectra were recorded with either L2MS or SIMS in multiple regions of the sample surface, to ensure the consistency of the method and to build a database allowing a more advanced statistical analysis. Once all the peaks coming from the substrate are removed, the data is ready to be processed.

### 3.1 Mass defect analysis

Mass defect analysis is used to assign a molecular formula to the recorded accurate mass<sup>50,51</sup>. By convention, the mass defect of <sup>12</sup>C is defined as zero, therefore the mass defect of every other existing isotope is either positive or negative, depending on its relative nuclear binding energy to <sup>12</sup>C. Since each nuclide has a unique mass defect, molecules with different isotopic composition have unique exact mass. For example, while a resolving power of around 5000 is sufficient to completely separate C<sub>14</sub>H<sub>10</sub><sup>+</sup> and C<sub>13</sub>H<sub>6</sub>O<sup>+</sup>, for closely spaced ions the required resolving power

can easily increase up to 10<sup>5</sup> or even higher. As the m/z increases, the number of combinations of different elements resulting in the same nominal mass grows very fast. This experimental limitation is already tackled in Irimiea *et al.*<sup>19</sup> when discussing the role of oxygen containing compounds. Nevertheless, a lower mass resolution mass spectrum can provide several helpful information. In particular, in the investigation of soot particles sampled from laboratory flames C, H and O are the major contributors to the total mass of soot, and therefore the mass analysis of peaks with a high signal-to-noise ratio (SNR) can be reasonably limited to C<sub>m</sub>H<sub>n</sub>O<sub>p</sub><sup>+</sup> ions. Identification within 5 ppm, often but not necessarily assumed as “certain”<sup>52</sup>, in our work is possible up to m/z ≈ 150 – 200. *A priori* knowledge of the samples and experimental conditions can extend this range up to m/z ≈ 500 – 550 and lead to self consistent results and coherence with many other works in the literature.

The mass defect analysis can also be used to simplify the visualisation of complex mass spectra (e.g. Figures S1 and S2). This is generally achieved by plotting the mass defects of all peaks versus their nominal mass. The resulting graph (mass defect plot, Figure 1 and S3) enables the visualisation of complex databases in one single plot, and highlights trends that are often invaluable to identify unknown species. For instance aliphatic, aromatic or

298 polycyclic aromatic hydrocarbons are aligned on different positive  
299 slopes corresponding to the addition of H atoms. When analysing  
300 samples containing hydrocarbons with different degrees of alky-  
301 lation, the Kendrick mass defect can be used as an alternative way  
302 of presenting the mass defect data<sup>50,51</sup>. Kendrick mass defect is  
303 calculated from the re-normalised mass of a repeating molecu-  
304 lar fragment to an integer value as shown in Equation 1 for the  
305 common case of CH<sub>2</sub> ( $m = 14.01565$ ):

$$m^{Kendrick} = m^{IUPAC} \frac{14.0000}{14.01565} \quad (1)$$

306 After this conversion, homologous series that contain the re-  
307 peating fragment have identical Kendrick mass defect and are  
308 found aligned on horizontal lines, making their identification  
309 even easier<sup>50,53</sup>. This is useful when dealing with repeating alkyl  
310 groups for instance, since their mass defect increases regularly  
311 with their molecular weight and makes their association to a cer-  
312 tain series less intuitive when represented on conventional mass  
313 defect plots<sup>50</sup>. The most convenient approach (conventional or  
314 Kendrick) heavily depends on the nature of the sample. If the  
315 sample is dominated by a variety of different species, the use of  
316 the conventional mass defect is more advisable. However, when  
317 the mass spectrum contains many species that only differ by a  
318 repeating unit such as aliphatic chains for instance (Table S1),  
319 Kendrick mass defect is more advantageous (Figure S4).

320 In this work, mass defect analysis is applied to the data ob-  
321 tained from L2MS and SIMS to demonstrate its effectiveness  
322 when dealing with a variety of mass spectrometric data. Figure  
323 1 shows the mass defect plot obtained from sample C2 analysed  
324 by L2MS. The suggested representation merges into one graph  
325 important information extracted from the raw mass spectra that  
326 include the peaks mass defect (y-axis), nominal mass (x-axis) and  
327 relative abundance (dot size). Species that line up in the mass  
328 defect plots typically contain a repeating unit. Additionally, the  
329 detection of a series of homologous species can help the identi-  
330 fication of unknown peaks. This is especially helpful for species  
331 with high molecular masses, where the attribution of a chemical  
332 formula can be rather delicate.

333 As PAHs exhibit a high thermodynamic stability<sup>54</sup>, they appear  
334 in great abundance in all mass spectra and this is amplified by the  
335 high sensitivity of the analysis technique to these specific com-  
336 pounds (Figure S1). Since the H/C ratio of PAHs is low com-  
337 pared to other hydrocarbons, they have a relatively small mass  
338 defect and are thus easily distinguishable from other hydrocar-  
339 bons. For instance, aromatic hydrocarbons that contain the same  
340 number of hydrogen atoms and progressively increasing number  
341 of carbon atoms (e.g. C<sub>10</sub>H<sub>8</sub> → C<sub>12</sub>H<sub>8</sub> → C<sub>14</sub>H<sub>8</sub> → ... → C<sub>22</sub>H<sub>8</sub>)  
342 can be found on the same horizontal line. Besides hydrocarbons,  
343 all samples contain oxygen and nitrogen organic derivatives to  
344 some extent. As a rule of thumb, in the mass defect plot of com-  
345 bustion generated aerosols, oxygen containing hydrocarbons are  
346 often found below the corresponding hydrocarbons due to the  
347 large negative mass defect of oxygen. Nitrogen containing hydro-  
348 carbons show distinct behaviours. For instance, organic amines  
349 are often found mixed to their corresponding hydrocarbons due  
350 to the nucleophilicity of nitrogen that results in their tendency to

bind one additional hydrogen atom post-ionisation. Organic ni-  
trates, on the other hand, tend to be found at lower mass defect  
due to the presence of oxygen.

Kendrick mass defect can be used to emphasise some less obvi-  
ous patterns as shown in Figure S4, in which CH ( $m = 13.007825$ )  
is used as the base unit.

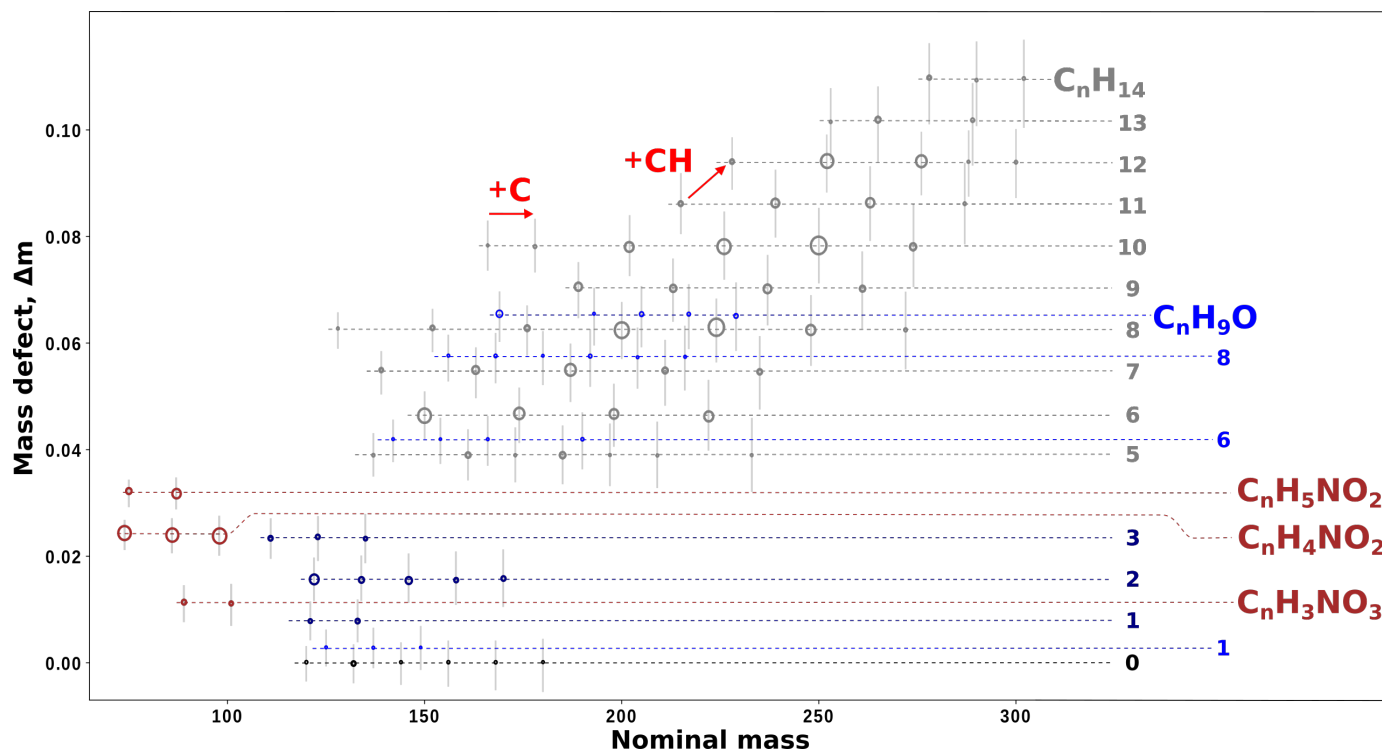
## 3.2 Statistical analysis

In this section we detail the chemometric techniques, based on  
commonly used statistical tools like multivariate analysis, that  
were adopted by our group to extract chemical information from  
mass spectrometric data. A mass spectrometry database can con-  
tain an extremely variable number of mass spectra (observations),  
and each of them typically contain up to thousands of peaks (vari-  
ables). This database structure should be taken into consideration  
when choosing the most appropriate statistical methods.

### 3.2.1 Principal component analysis

PCA is a powerful statistical tool that can be used to classify sam-  
ples and reveal trends and patterns in databases<sup>5</sup>, and is often  
used to increase the readability of very complex data<sup>55</sup>. PCA ap-  
plied to mass spectrometry is especially useful when many mass  
spectra are being compared, since it reduces the dimensionality  
of the database while preserving most of the original informa-  
tion. PCA is a non-parametric analysis, i.e. its output is inde-  
pendent of any hypotheses about data distribution<sup>56</sup>. In this  
work, PCA is performed on a matrix containing the integrated  
peaks (variables) against the samples (observations). Before ap-  
plying PCA, data obtained from mass spectrometry should un-  
dergo a special preparation procedure<sup>56,57</sup> that includes calibra-  
tion, baseline removal, construction of a peak list, peak integra-  
tion and standardisation. PCA applied to data with no normali-  
sation/standardisation is mostly affected by the largest raw vari-  
ance, which can skew the overall interpretation of the dataset.  
Therefore, normalisation techniques are applied to mass spec-  
tra prior PCA analysis when there are differences in the sam-  
ples weight, volume or other properties that may result in ad-  
ditional sources of variance. The most popular and generally rec-  
ommended normalisation method is the normalisation to the total  
ion count (TIC), i.e. the integrated ion count over a given mass  
range<sup>18,58,59</sup>.

Care has to be taken when building the peak list as it should  
only contain species representative of the sample. Minor-  
abundance isotopes are usually excluded from the peak list, thus  
allowing to focus on the major-abundance isotopic species<sup>58</sup>.  
Peaks coming from the substrate and/or originate from the  
sample-substrate reactivity should also be disregarded. Identi-  
fying these peaks, especially the ones corresponding to reaction  
products, can be a difficult task. One approach to their identifica-  
tion involves comparing mass spectra of the sample deposited in  
the same experimental conditions but on different substrates (e.g.  
Si and Ti wafers)<sup>18</sup>. Another possibility relies on the use of PCA:  
species coming from the sample-substrate reactivity become less  
prominent as the substrate coverage increases and is less avail-  
able for the reaction, and are thus likely to be found all clustered  
in the same principal component.



**Fig. 1** Mass defect plot obtained from the L2MS mass spectrum of miniCAST soot, C2 sample. The data points represent the assigned accurate mass. The size of the data points is proportional to the corresponding peak integrated area, normalised to the total ion count after background subtraction. Molecular formulas of homologous species are displayed. The error bars show the uncertainty on the accurate mass calculated from the obtained mass resolution.

Each principal component (PC) accounts for a defined percent-<sup>432</sup> age of the total variance within the data set, are represented in a<sup>433</sup> scree plot and used to select the PCs to take into consideration. The loadings represent the weights of each variable used to cal-<sup>434</sup> culate the PCs, and are used to understand the contribution of each variable to the selected PC. The distance of an observation<sup>435</sup> from a PC is represented on the scores plot. Scores are obtained<sup>436</sup> for each observation in the database and for each principal com-<sup>437</sup> ponent, and are often used as a base to display and classify the<sup>438</sup> samples. In the score plot, similar observations group together<sup>439</sup> and are separated from dissimilar observations. The clustering<sup>440</sup> of the scores is strongly related to the values of the loadings, and<sup>441</sup> they are discussed as a whole. The most challenging part of PCA is<sup>442</sup> the interpretation of individual PCs and their contribution to the<sup>443</sup> investigated processes. To this purpose, there is a vast literature<sup>444</sup> providing general guidelines that should be followed<sup>5,60-62</sup>.<sup>445</sup>

To illustrate the potential of this technique, we show below<sup>446</sup> some application to mass spectrometric data of various combus-<sup>447</sup> tion generated aerosol samples.<sup>448</sup>

### 3.2.1.1 MiniCAST soot, L2MS

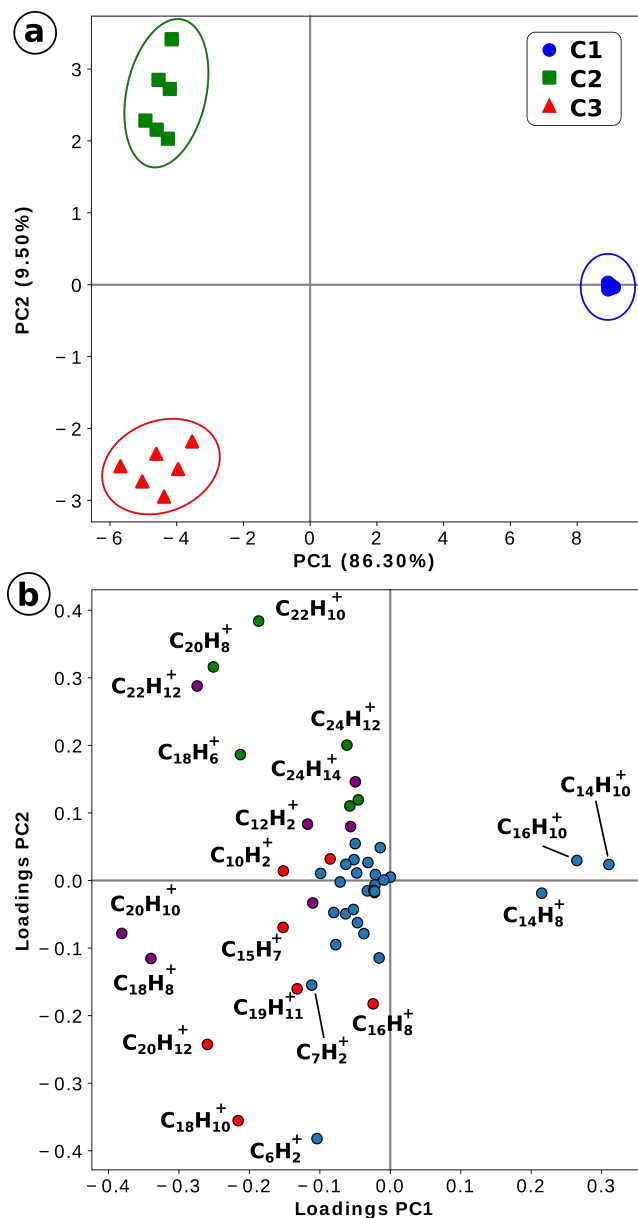
When L2MS mass spectra of miniCAST soot samples are exam-<sup>451</sup> ined, PC1 and PC2 account for  $\sim 96\%$  of the total variance, and<sup>452</sup> are therefore only considered for the data interpretation. The<sup>453</sup> three samples are well separated in the PC2 vs. PC1 scores plot<sup>454</sup> (Figure 2). Sample C1 is highly influenced by  $C_{14}H_8$ ,  $C_{14}H_{10}$  and<sup>455</sup>  $C_{16}H_{10}$  (high positive PC1 scores) whereas C2 and C3 are domi-<sup>456</sup> nated by higher mass aromatic compounds (negative PC1 scores).<sup>457</sup>

It can be noticed that PC2 ( $\sim 10\%$ ) allows for better discrimina-  
tion between the samples than PC1, especially C2 and C3.

### 3.2.1.2 Flame and ICE soot, SIMS

PCA is applied to the ensemble of SIMS mass spectra obtained  
in positive polarity from soot samples generated by the gasoline  
engine and the laboratory flame (diesel and kerosene fuels). PC1  
and PC2 account together for the 73.3% of the total variance.  
Two main groups are observed in the score plot of both positive  
and negative ions (Figure 3 and S5). While it was not possible to  
clearly associate a phenomenon to PC1 (51.7% of total variance),  
the samples are well separated by the different emission source  
(engine, GOM, and flame, D and K) in PC2 (21.6% of total vari-  
ance). At this level of the analysis PCA cannot distinguish soot  
generated by burning the two different liquid fuels (diesel and  
kerosene) in laboratory flames, which appear mixed together in  
negative PC2.

PC1 is mainly associated to high H/C fragment ions (negative  
contribution, red dots in the loadings plot (Figure 3), and low  
H/C fragment ions probably resulting from the dissociation of  
large aromatic hydrocarbons (positive contribution, green dots  
in the loadings plot). The main contributions to PC2 come from  
aromatic species (positive contribution, blue dots on the loadings  
plot), and to a smaller extent to high H/C fragment ions. There-  
fore, the contribution of high H/C fragment ions, possibly related  
to the dissociation of aliphatic hydrocarbons, depends less on the  
fuel and more on the combustion conditions (engine vs. con-



**Fig. 2** Score plots of PC2 vs PC1 for miniCAST soot samples obtained with L2MS – (a). Ellipses highlight data points coming from different samples and are added for visual purposes only. (b) – the corresponding loadings plot of PC2 vs PC1. Several homologous series are highlighted:  $C_{n+8}H_n$  – red,  $C_{n+10}H_n$  – purple,  $C_{n+12}H_n$  – green.

controlled laboratory flames).

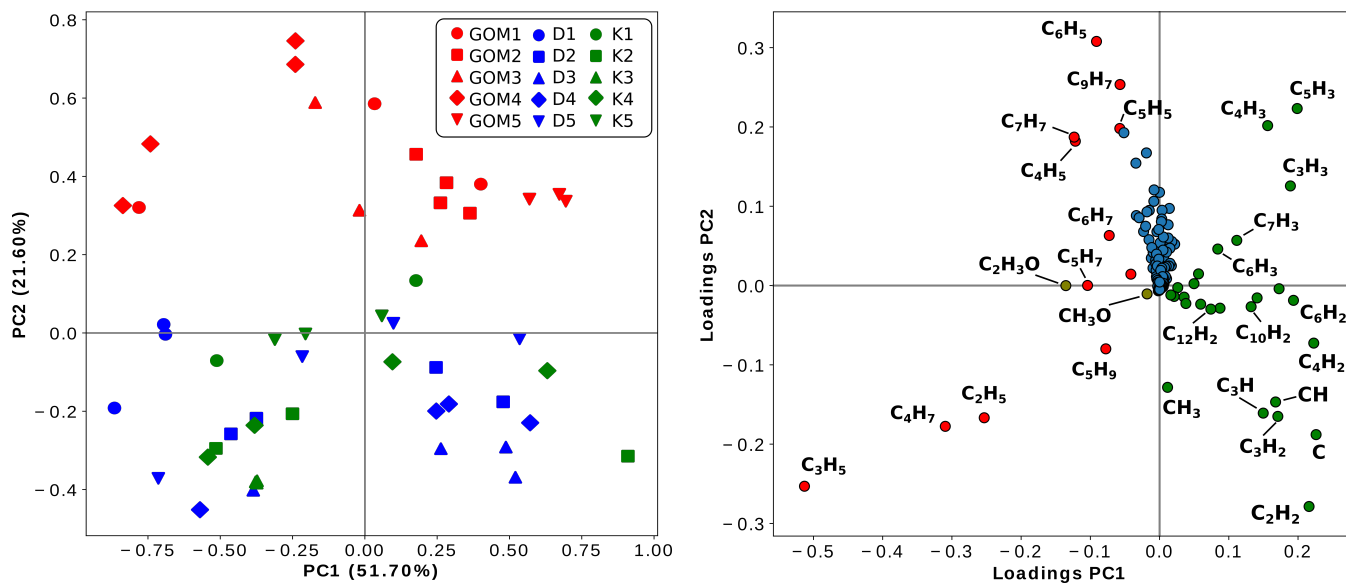
Going a step further, PCA is applied to gasoline soot samples obtained in different engine regimes in order to determine their impact on the chemical composition (Figure 4). There is an obvious separation between normal engine operation regimes (GOM, GOH) and the ones which simulate a malfunction (GEF, GEO). A good discrimination is achieved with only the first two components that account for  $\sim 98\%$  of the total variance. PC1 alone ( $\sim 91\%$ ) allows the separation of regimes, based on the abundance of aliphatic fragment ions (positive contribution to PC1, marked in red in Figure S6). Consequently, samples that simulate a malfunction (GEF, GEO) are characterised by a higher relative

contribution from aliphatic fragment ions compared to optimised engine regimes (GOM, GOH). PC2 is linked to the contribution of aliphatic fragment ions and aromatic species (positive PC2 value), however some aliphatic fragment ions ( $C_5H_7$ ,  $C_5H_9$ ,  $C_3H_7$ ,  $C_4H_7$ ) show a contribution to negative PC2). The data points corresponding to optimal engine regimes form a smaller cluster. This implies that soot produced in conditions simulating engine malfunction shows a much larger variability in chemical composition.

At this point of the analysis, it is clear that the two regimes that simulate a malfunction (GEF, GEO) exhibit similarities, while being well separated from the optimised regimes (upper panel of Figure 4). This implies that the variance of a certain principal component for them is much smaller than the one responsible for the separation between optimised and non-optimised regimes. Consequently, each group should be analysed independently, thus uncovering even smaller contributions to the variance. To demonstrate this concept, the same statistical method was applied a second time to the two non-optimised regimes, and their comparison lead to discriminate between the two main contributors to particulate emissions of the internal combustion engine: fuel and oil, Figure 4. In this case, PC1 ( $\sim 71\%$ ), accountable for the separation of the two regimes, is linked to the contribution of hydrogen-rich hydrocarbons on one side (negative contribution) and of fragment ions and aromatic species on the other (positive contribution). This reveals that oil-related soot particles feature more hydrogen-rich hydrocarbons, while an excess of gasoline leads to the production of more aromatic species, Figure S6. The increase of the contribution of fragment ions in the latter is probably linked to the increase in the aromatic contribution, since the majority of fragment ions can be related to dissociation reactions of PAHs<sup>63</sup>. PC2 ( $\sim 20\%$ ) is associated to the presence of aromatic hydrocarbons (blue dots in Figure S6). One can also notice that samples corresponding to the engine regime with a low air/fuel ratio (GEF1) surprisingly lie in the oil-excess region, while samples GEO3 appear far from the oil-excess region (Figure 4). It is likely that the specific behaviour observed for these samples relates to their particle size (Table 1) but correlating size to chemical composition is out of scope of this paper and will be addressed in a future work.

### 3.2.2 Hierarchical clustering analysis

Hierarchical clustering analysis (HCA) is a MVA method that identifies patterns in a dataset by creating groups of observations called clusters. Unlike PCA, HCA accounts for the total variance in the database<sup>60,62</sup>. HCA is based on a simple approach for building the clusters that starts with one cluster for each observation and finishes with a single cluster containing the entire database. At each step, the two closest clusters are merged into a single new cluster resulting in a dendrogram representative of the database. In order to decide which clusters to merge, different approaches to measure their distance can be used and give rise to several hierarchical methods<sup>61,62</sup>. In this work, HCA (group average method, Euclidean distances) is applied to the same standardised matrix used for PCA analysis, on both columns (observations) and rows (variables). The HCA output is built in a heatmap organised by the clusters obtained on observations and variables.



**Fig. 3** Score plot of PC2 vs PC1 for positive ions of soot samples obtained from gasoline engine and laboratory flames (left panel). Corresponding loadings plot of the first two principal components (right panel). For sample description see Table 1.

525 This representation improves the visualisation of clusters in the 559  
 526 multidimensional space, in which each tile represents the value 560  
 527 of the correlation between observations and variables. 561

528 The heatmap obtained for the samples analysed in SIMS pos- 562  
 529 itive polarity is shown in Figure 5. HCA groups the samples in 563  
 530 three main clusters (C1, C2 and C3) at distance d1 function of 564  
 531 the characteristics of the five clusters of variables (R1, R2, R3, R4 565  
 532 and R5). Cluster C1 is specific to samples GEO1-4, GOM4 and 566  
 533 D1 due to the high contribution of compounds with  $H/C > 1$  and 567  
 534 identified in the C1-1 cluster. C1-2 is dissimilar from the C1-1 due 568  
 535 to the presence of aromatic hydrocarbons and other compounds 569  
 536 with low  $H/C$  ratio. Soot collected from the gasoline engine in 570  
 537 optimal conditions and after the addition of oil are dominated by 571  
 538 R5, while there is a shift to R1 and R2 for soot collected from the 572  
 539 diesel flame. Contrary to C1, C2 has a high contribution of frag- 573  
 540 ment ions with high (R4) and low (R1)  $H/C$  ratio. C2 shows that 574  
 541 soot collected from the engine in optimal conditions with high 575  
 542 and medium load have similar chemical fingerprint. 576

543 This representation offers at once a clustering of the samples 577  
 544 function of the three main classes of chemical compounds iden- 578  
 545 tified in the mass spectra. For instance, the high content of aro- 579  
 546 matic hydrocarbons and low  $H/C$  fragment ions is specific to soot 580  
 547 collected from the kerosene flame. Basically, the addition of oil 581  
 548 increases the fraction of high  $H/C$  fragment ions in the emissions, 582  
 549 the normal operation conditions of the engine have an intermedi- 583  
 550 ate content of high  $H/C$  fragment ions and a slight contribution of 584  
 551 aromatics with four and five aromatic rings, while kerosene soot 585  
 552 contains the highest contribution of aromatic compounds and low 586  
 553  $H/C$  fragment ions. HCA is also applied to L2MS and SIMS neg- 587  
 554 ative polarity data as detailed in the Supplementary Information. 588  
 555 In this work, HCA is applied to the raw data corresponding to the 589  
 556 selected mass spectra but its usefulness can be extended to more 590  
 557 compact data after using another statistical method for sorting 591  
 558 the input variables and observations. One of the advantages of 592

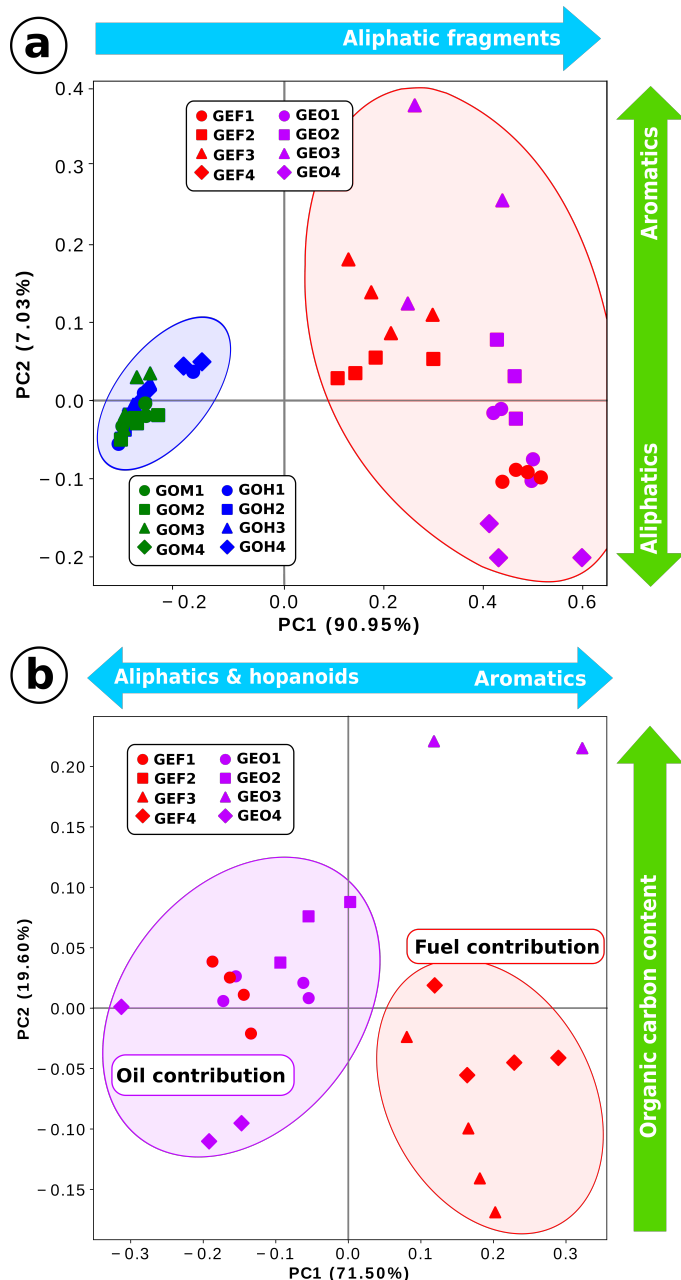
this method is that it does not require the raw data set. Moreover, HCA can be used to visualise clusters that form in the principal component space, after applying the PCA, or it can group samples according to other properties (mass defect, contribution from different classes of compounds, etc).

### 3.3 Mass peaks grouping into chemical classes

A detailed description of the soot chemical composition is certainly desirable and can lead to important clues on the soot formation, growth, ageing and reactivity. However, this can rapidly turn into a very cumbersome task, especially if many different samples are analysed. For the sake of simplicity, most of the time, and especially when long time-series of field-collected data are to be treated, individual compounds are grouped in classes (e.g. aliphatics, aromatics, oxygenated, sulphur-containing hydrocarbons and so on). This grouping of mass peaks into appropriate classes allows easier comparison with other experimental measurements (e.g. OC/EC<sup>29</sup>) and facilitates the interaction with modellers that use the data as inputs for various scales simulations. Moreover, this grouping of peaks is also useful when mass spectra of several samples are compared to each other in order to reveal general trends in their chemical composition.

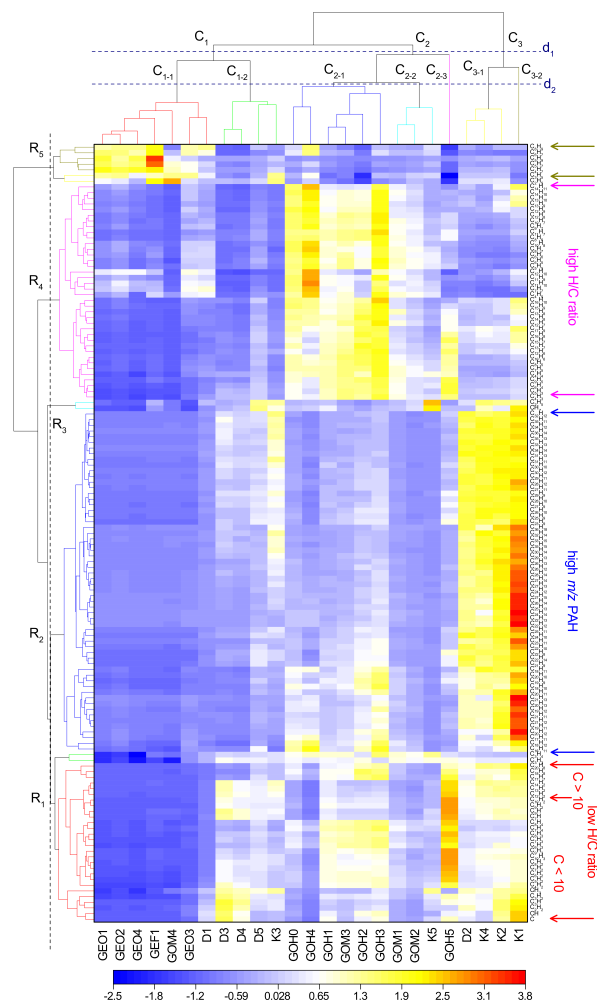
When it comes to the chemical composition of combustion generated aerosols, three non-specific indicators are often considered: amount of ash components (inorganic compounds, IC), amount of carbon associated to the carbonaceous matrix (elemental carbon, EC), and amount of carbon found in organic compounds (organic carbon, OC)<sup>64</sup>. IC alone can sometimes help identify the main source of the emissions. For instance,  $K^+$ ,  $Na^+$ ,  $K_2Cl^+$  and  $K_3SO_4^+$  in the positive polarity mass spectra and  $Cl^-$ ,  $SO_3^-$ ,  $HSO_4^-$  and  $KCl_2^-$  in the negative polarity mass spectra are known to be markers of wood combustion<sup>65</sup>. Generally speaking, since IC potentially contains many inorganic compounds, it can





**Fig. 4** Score plots of the first two principal components for soot samples produced by a single cylinder engine. Upper panel – discrimination between different engine regimes, lower panel – particle source discrimination. Ellipses highlight clusters of data points and are for visual purposes only. For sample description see Table 1.

and should be further broken down into source specific groups when characterising complex systems such as internal combustion engines. In this case, accepted grouping of inorganic compounds is: fuel specific (compounds that are coming from fuel additives and trace elements (Na, K)<sup>11,66</sup>), oil specific (detergent and anti-wear additives (P, Ca)<sup>67</sup> and engine wear tracers (Fe, Al, Cr)<sup>30,67,68</sup>). For addressing the elemental carbon (EC) component, carbon clusters  $C_n^-$  ( $n=2-4$ ) are considered to be appropriate markers in aerosol mass spectrometry<sup>64</sup>. This is also confirmed by the high positive correlation between  $C_2^-$ ,  $C_3^-$  and  $C_4^-$



**Fig. 5** Two-way hierarchical clustering heat-map for positive ions of gasoline, diesel and kerosene soot obtained with SIMS. Each column corresponds to the averaged mass spectra obtained for a soot sample. The contribution of each mass in individual samples is expressed as relative value and is represented by the cell colour.

signals in the recorded mass spectra<sup>26</sup>. In single particle mass spectrometry, carbon clusters with even higher masses are also considered to be representative of the elemental carbon ( $C_5^-$  at 60 u,  $C_6^-$  at 72 u and  $C_7^-$  at 84 u)<sup>11</sup>. While the handling of IC and EC is relatively straightforward, the OC landscape looks far more complex, with an overwhelming variety of organic compounds, generated in various processes and being themselves main actors of broad-range time-scale reactivity. A subsequent classification of different organic species according to their functional group(s) seems therefore necessary. However, the detailed chemical analysis of a complex mixture of chemicals based on mass spectrometric data only is still an important challenge that requires the identification of the individual ion dissociation patterns. On a practical ground, being able to distinguish these compounds is very important since they all have different sources and roles in the soot formation and ageing mechanisms. For instance, PAHs form during combustion and are well known as building blocks of soot particles and are generally seen as reliable markers of the over-

619 all OC content<sup>29</sup>. Organic hydroxyl groups are linked to alcohols  
 620 that are commonly used as additives in gasoline. The presence of  
 621 many compounds containing carbonyl groups has been proposed  
 622 as a marker to distinguish fresh emissions from soot particles aged  
 623 in the atmosphere<sup>69</sup>.

624 A combination of previously described mass peak classification  
 625 methods is shown in Table 2 along with chemical formula assign-  
 626 ments<sup>63</sup>. Detailed classification of molecular ions by functional  
 627 groups remains difficult by MS alone, however it can be achieved  
 628 in combination with complementary techniques (e.g. FTIR).<sup>26</sup>  
 629 Also, for the sake of simplicity, Table 2 displays only the nominal  
 630 masses, but the peak assignment is based on the exact mass (see  
 631 mass defect analysis, Section 3.1). The discussion below is based  
 632 on this grouping of mass peaks.

633 Depending on the studied samples, the analysis will focus on  
 634 specific classes from Table 2. For soot samples obtained with  
 635 the miniCAST standard generator, one may want to address the  
 636 impact of the oxidation flow. A possible focus is therefore on  
 637 the evolution of the oxygenated species vs. PAHs (linked to the  
 638 OC content). Since miniCAST soot is a well-studied standard, it  
 639 also allows the comparison of mass spectrometric results with the  
 640 ones reported in the literature based on other experimental tech-  
 641 niques. In the present case, Figure 6 clearly shows an increase of  
 642 the oxygenated species abundance with the oxidation flow, how-  
 643 ever a low oxidation flow (C2 and C3) leads to the formation of  
 644 more PAHs, which confirms previous observations on the same  
 645 set-points of the miniCAST generator<sup>43,70</sup>.

646 Even though examining trends for specific groups can be very  
 647 informative, when it comes to complex mass spectra containing a  
 648 multitude of peaks that can be separated in many different ways,  
 649 not all the groups feature useful trends. It is therefore advis-  
 650 able to first identify the species of interest, groups or individual  
 651 compounds that can be linked to variations in the chemical com-  
 652 position of the samples. This information can be retrieved from  
 653 PCA and HCA as discussed in the sections 3.2.1 and 3.2.2, re-  
 654 spectively. Based on the statistical analysis of positive polarity  
 655 SIMS mass spectra of gasoline, diesel and kerosene soot samples,  
 656 three groups of interest are chosen for further analysis as shown  
 657 in Figure 6: low-mass and low H/C ions (from the dissociation  
 658 of aromatic species<sup>63</sup>), low-mass and high H/C ions (from the  
 659 dissociation of aliphatic species), and finally large aromatic ions  
 660 (mostly PAHs, stable enough to be detected as molecular ions).  
 661 Gasoline soot shows higher content of large aromatic compounds,  
 662 with high and almost constant contribution to all considered par-  
 663 ticle sizes. Gasoline soot also features the least fragmentation  
 664 that is well consistent with the higher contribution of large ar-  
 665 omatics if compared to diesel and kerosene soot. For the other  
 666 two fuels, different zones of the flame, corresponding to different  
 667 stages in the soot formation process, were probed, therefore the  
 668 variation in aromatic content looks more pronounced. It is clear  
 669 that the aliphatic content alone cannot be used to discriminate<sup>675</sup>  
 670 between soot coming from combustion of different fuels, just like<sup>676</sup>  
 671 it was concluded from PCA. However, it still provides valuable in-<sup>677</sup>  
 672 formation about different soot maturity. For example, for diesel<sup>678</sup>  
 673 soot the contribution of aliphatics gradually increases with the<sup>679</sup>  
 674 sampled HAB ( $HAB \geq 12\text{ cm}$ ). On the other hand the HCA on<sup>680</sup>

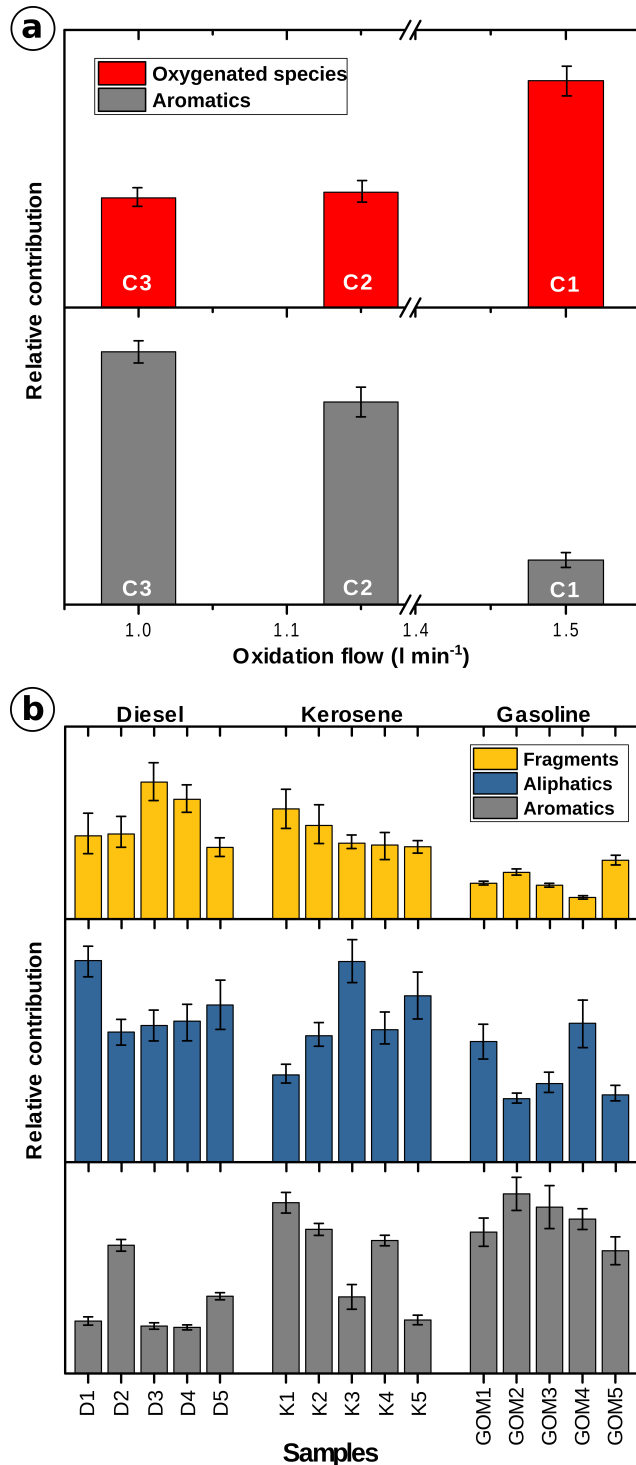


Fig. 6 Several trends retrieved from mass spectra of: (a) – miniCAST soot (L2MS), (b) – gasoline, diesel and kerosene soot (SIMS).

the negative polarity of SIMS is much easier to interpret because the results clearly discriminate the laboratory flame soot from the one produced with the gasoline engine. The samples belonging to the latest category are clearly evidenced by the presence of sulphur and oxygen containing compounds while the soot from the flames contains mainly OC and EC. Generally speaking, the

**Table 2** Grouping of mass peaks into chemical classes

Category	m/z	Formula	m/z	Formula	m/z	Formula	m/z	Formula
Aliphatics (alkynes, alkene, alkyl, etc.)	15	$CH_3$	54	$C_4H_6$	71	$C_5H_{11}$	99	$C_7H_{15}$
	27	$C_2H_3$	55	$C_4H_7$	81	$C_6H_9$	109	$C_8H_{13}$
	29	$C_2H_5$	57	$C_4H_9$	83	$C_6H_{11}$	111	$C_8H_{15}$
	41	$C_3H_5$	67	$C_5H_7$	85	$C_6H_{13}$	113	$C_8H_{17}$
	43	$C_3H_7$	68	$C_4H_8$	95	$C_7H_{11}$		
	53	$C_4H_5$	69	$C_5H_9$	97	$C_7H_{13}$		
Aromatics	26	$C_2H_2$	64	$C_5H_4$	152	$C_{12}H_8$	216	$C_{17}H_{12}$
	38	$C_3H_2$	74	$C_6H_2$	154	$C_{12}H_{10}$	228	$C_{18}H_{12}$
	39	$C_3H_3$	75	$C_6H_3$	166	$C_{13}H_{10}$	252	$C_{20}H_{12}$
	40	$C_3H_4$	76	$C_6H_4$	178	$C_{14}H_{10}$	276	$C_{22}H_{12}$
	50	$C_4H_2$	78	$C_6H_6$	266	$C_{21}H_{14}$	278	$C_{22}H_{14}$
	51	$C_4H_3$	91	$C_7H_7$	190	$C_{15}H_{10}$		
	63	$C_5H_3$	128	$C_{10}H_8$	202	$C_{16}H_{10}$		
O-containing (carbonyls, acids, ethers, alcohols, etc.)	31	$CH_3O$	69	$C_4H_5O$	87	$C_5H_{11}O$	129	$C_7H_{13}O_2$
	33	$CH_5O$	71	$C_4H_7O$	89	$C_5H_{13}O$	137	$C_{10}HO$
	43	$C_2H_3O$	73	$C_3H_5O_2$	97	$C_6H_9O$	142	$C_{10}H_6O$
	45	$C_2H_5O$	73	$C_4H_9O$	97	$C_5H_5O_2$	156	$C_{11}H_8O$
	47	$CH_3O_2$	75	$C_3H_7O_2$	101	$C_6H_{13}O$	166	$C_{12}H_6O$
	47	$C_2H_7O$	75	$C_4H_{11}O$	105	$C_7H_5O$	169	$C_{11}H_9O$
	53	$C_4H_5$	81	$C_5H_5O$	109	$C_7H_9O$	180	$C_{13}H_8O$
	55	$C_3H_3O$	83	$C_5H_7O$	111	$C_6H_7O_2$	205	$C_{14}H_9O$
	57	$C_3H_5O$	85	$C_5H_9O$	111	$C_7H_{11}O$		
	59	$C_3H_7O$	85	$C_4H_5O_2$	119	$C_8H_7O$		
	61	$C_2H_5O_2$	87	$C_5H_{11}O$	123	$C_7H_7O_2$		
	61	$C_3H_9O$	87	$C_4H_7O_2$	125	$C_9HO$		
N-containing	26	$CN$	46	$CH_4NO$	60	$C_2H_6NO$	89	$C_2H_3NO_3$
	29	$CH_3N$	55	$C_3H_5N$	74	$C_2H_4NO_2$	98	$C_4H_4NO_2$
	44	$CH_2NO$	55	$C_2H_3N_2$	87	$C_3H_5NO_2$	121	$C_8H_{11}N$
S-containing	32	$S$	44	$CS$	46	$CH_2S$		
Unclassified hydrocarbons	28	$C_2H_4$	56	$C_4H_8$	84	$C_6H_{12}$	112	$C_8H_{16}$
	42	$C_3H_6$	70	$C_5H_{10}$	98	$C_7H_{14}$		

trends that are shown herewith are very useful when interpreting the data. However, they are almost impossible to notice in the raw mass spectra. Being able to follow the contribution of a group of related molecules hidden in a much larger ensemble of signals is a powerful feature used to uncover trends that would have remained hidden to a more basic analysis. The fact that PCA and HCA are able to separate the selected samples into categories dependent on their unique pattern of chemical signatures proves that mass spectrometry and MVA provide useful insights into their properties. The usefulness of this approach allows for an easier identification and traceability of combustion generated particles with unknown sources.

## 4 Conclusions

Our recently developed comprehensive methodology (based on mass defect analysis, PCA/HCA multivariate methods) dedicated to the chemical analysis of combustion-generated aerosols is applied here to the study of 30 soot samples generated by three different sources using four different fuels. Laser and secondary ion mass spectrometry techniques are used to probe their surface chemistry. A few examples on the performances of this methodology are provided, showcasing its ability to clearly discriminate samples according to various parameters, such as com-

bustion source, soot maturity, or engine operating conditions. The correlations evidenced by the MVA methods were used for peak clustering to highlight the evolution of grand chemical classes with the combustion conditions. These trends, along with detailed molecular-level information, can further help constrain the processes involved in particulate matter emissions and predict the impact of soot particles on the environment and human health. Moreover, aiming for a standardised (generally accepted) methodology in treating complex mass spectrometry data in aerosol science would certainly allow easier intercomparison and the building of extensive shared databases for further specific developments. An appealing perspective is the possible application of neural networks to this type of big data, which would lead to great advances in automated real-time processing of large dataflows.

## Conflicts of interest

There are no conflicts to declare.

## Acknowledgements

This work was supported by the French National Research Agency (ANR) through the PIA (Programme d'Investissement d'Avenir) under contract ANR-10-LABX-005 (CaPPA – Chemical and Phys-

ical Properties of the Atmosphere), the European Commission  
Horizon 2020 project PEMS4Nano (H2020 Grant Agreement  
#724145), and the CLIMBIO project via the Contrat de Plan Etat  
Région of the Hauts-de-France region. We thank N. Nuns from the  
Regional Surface Analysis Platform for assistance with the SIMS  
measurements.

## References

- 1 R. S. Plumb *et al.*, *Rapid Communications in Mass Spectrometry*, 2003, **17**, 2632–2638.
- 2 M. Brokl *et al.*, *Journal of Chromatography A*, 2014, **1370**, 216–229.
- 3 S. Shrestha and F. Kazama, *Environmental Modelling and Software*, 2007, **22**, 464–475.
- 4 W. K. Härdle and L. Simar, *Applied Multivariate Statistical Analysis Course*, Springer, 2015.
- 5 H. Abdi and L. J. Williams, *English*, 2010, **2**, 433–470.
- 6 J. L. Jimenez *et al.*, *Science*, 2009, **326**, 1525–1529.
- 7 C. Fountoukis *et al.*, *Atmospheric Chemistry and Physics*, 2014, **14**, 9061–9076.
- 8 M. Crippa *et al.*, *Atmospheric Chemistry and Physics*, 2014, **14**, 6159–6176.
- 9 T. B. Onasch *et al.*, *Aerosol Science and Technology*, 2012, **46**, 804–817.
- 10 A. Faccinnetto *et al.*, *Combustion and Flame*, 2011, **158**, 227–239.
- 11 U. Kirchner *et al.*, *Journal of Aerosol Science*, 2003, **34**, 1323–1346.
- 12 N. Mayama *et al.*, *Analytical Sciences*, 2013, **29**, 479–482.
- 13 F. Aubriet and V. Carré, *Analytica Chimica Acta*, 2010, **659**, 34–54.
- 14 S. Eliuk and A. Makarov, *Annual Review of Analytical Chemistry*, 2015, **8**, 61–80.
- 15 K. Wang *et al.*, *Atmospheric Environment*, 2018, **189**, 22–29.
- 16 C. Zuth *et al.*, *Analytical Chemistry*, 2018, **90**, 8816–8823.
- 17 J. Cain *et al.*, *Physical Chemistry Chemical Physics*, 2014, **16**, 25862–25875.
- 18 C. Irimiea *et al.*, *Rapid Communications in Mass Spectrometry*, 2018, **32**, 1015–1025.
- 19 C. Irimiea *et al.*, *Carbon*, 2019, **144**, 815–830.
- 20 M. F. Heringa *et al.*, *Atmospheric Chemistry and Physics*, 2011, **11**, 5945–5957.
- 21 J. L. Jimenez, *Journal of Geophysical Research*, 2003, **108**, 8425.
- 22 M. T. Timko *et al.*, *Journal of Engineering for Gas Turbines and Power*, 2010, **132**, 061505.
- 23 W. B. Knighton *et al.*, *Journal of Propulsion and Power*, 2007, **23**, 949–958.
- 24 M. T. Timko *et al.*, *Combustion Science and Technology*, 2011, **183**, 1039–1068.
- 25 M. Bente *et al.*, *Analytical Chemistry*, 2008, **80**, 8991–9004.
- 26 S. Gilardoni *et al.*, *Journal of Geophysical Research Atmospheres*, 2017, **34**, 401–409.
- 27 A. Kortelainen, *PhD thesis*, University of Eastern Finland, 2016.
- 28 M. Abegglen *et al.*, *Atmospheric Environment*, 2016, **134**, 181–197.
- 29 D. Delhay *et al.*, *Journal of Aerosol Science*, 2017, **105**, 48–63.
- 30 J. Moldanová *et al.*, *Atmospheric Environment*, 2009, **43**, 38–44.
- 31 C. Giorio *et al.*, *Atmospheric Environment*, 2012, **61**, 316–326.
- 32 S. S. Lim *et al.*, *The Lancet*, 2012, **380**, 2224–2260.
- 33 V. Samburova, B. Zielinska and A. Khlystov, *Toxics*, 2017, **5**, 17.
- 34 R. Niranjana and A. K. Thakur, *Frontiers in Immunology*, 2017, **8**, 1–20.
- 35 T. Petry, P. Schmid and C. Schlatter, *Chemosphere*, 1996, **32**, 639–648.
- 36 D. A. Knopf, P. A. Alpert and B. Wang, *ACS Earth and Space Chemistry*, 2018, **2**, 168–202.
- 37 A. Faccinnetto *et al.*, *Environmental Science and Technology*, 2015, **49**, 10510–10520.
- 38 R. Lemaire *et al.*, *Proceedings of the Combustion Institute*, 2009, **32**, 737–744.
- 39 Y. Bouvier *et al.*, *Proceedings of the Combustion Institute*, 2007, **31 I**, 841–849.
- 40 P. Parent *et al.*, *Carbon*, 2016, **101**, 86–100.
- 41 F. X. Ouf *et al.*, *Scientific Reports*, 2016, **6**, 1–12.
- 42 J. Yon *et al.*, *Combustion and Flame*, 2018, **190**, 441–453.
- 43 R. H. Moore *et al.*, *Aerosol Science and Technology*, 2014, **48**, 467–479.
- 44 R. Zimmermann *et al.*, *Environmental Science and Technology*, 2001, **35**, 1019–1030.
- 45 O. P. Haefliger and R. Zenobi, *Analytical chemistry*, 1998, **70**, 2660–2665.
- 46 K. Thomson *et al.*, *Applied Surface Science*, 2007, **253**, 6435–6441.
- 47 A. Faccinnetto *et al.*, *Applied Physics A: Materials Science and Processing*, 2008, **92**, 969–974.
- 48 C. Mihean *et al.*, *Chemical Physics Letters*, 2006, **423**, 407–412.
- 49 C. Mihean *et al.*, *Journal of Physics: Condensed Matter*, 2008, **20**, 25221.
- 50 L. Sleno, *Journal of Mass Spectrometry*, 2012, **47**, 226–236.
- 51 C. A. Hughey *et al.*, *Analytical Chemistry*, 2001, **73**, 4676–4681.
- 52 A. G. Brenton and A. R. Godfrey, *Journal of the American Society for Mass Spectrometry*, 2010, **21**, 1821–1835.
- 53 R. Hilbig and R. Wallenstein, *Applied optics*, 1982, **21**, 913–917.
- 54 S. E. Stein and A. Fahr, *Journal of Physical Chemistry*, 1985, **89**, 3714–3725.
- 55 T. Adam, R. R. Baker and R. Zimmermann, *Journal of Agricultural and Food Chemistry*, 2007, **55**, 2055–2061.
- 56 Y. Tanaka, *Communications in Statistics - Theory and Methods*, 1988, 37–41.

- 830 57 R. E. Peterson and B. J. Tyler, *Atmospheric Environment*, 2002,  
831 36, 6041–6049.
- 832 58 P. Cejnar *et al.*, *Rapid Communications in Mass Spectrometry*,  
833 2018, 32, 871–881.
- 834 59 T. Alexandrov, *BMC Bioinformatics*, 2012, 13, S11.
- 835 60 L. Pei *et al.*, *Energy and Fuels*, 2008, 22, 1059–1072.
- 836 61 P. Reitz *et al.*, *Journal of Aerosol Science*, 2016, 98, 1–14.
- 837 62 R. Alvin C, *Methods of multivariate analysis - Second Edition*,  
838 Wiley - Interscience, 2001, pp. 1–727.
- 839 63 F. W. McLafferty and F. Tureek, *Interpretation of Mass Spectra*,  
840 University Science Books, Mill Valley, CA, 1993.
- 841 64 J. Pagels *et al.*, *Journal of Geophysical Research: Atmospheres*,  
842 2013, 118, 859–870.
- 843 65 J. Pagels *et al.*, *Journal of Aerosol Science*, 2003, 34, 1043–  
844 1059.
- 845 66 T. R. Dallmann *et al.*, *Atmospheric Chemistry and Physics*,  
846 2014, 14, 7585–7599.
- 847 67 E. S. Cross *et al.*, *Journal of Engineering for Gas Turbines and*  
848 *Power*, 2012, 134, 72801.
- 849 68 K. Aras, *Atmospheric Environment*, 1994, 28, 1385–1391.
- 850 69 S. Gilardoni *et al.*, *Journal of Geophysical Research: Atmo-*  
851 *spheres*, 2007, 112, 1–11.
- 852 70 J. Yon, A. Bescond and F.-X. Ouf, *Journal of Aerosol Science*,  
853 2015, 87, 28–37.