

1 **Automated compound speciation, cluster analysis, and**  
2 **quantification of organic vapors and aerosols using**  
3 **comprehensive two-dimensional gas chromatography and**  
4 **mass spectrometry**

5 Xiao He<sup>1</sup>, Xuan Zheng<sup>1\*</sup>, Shuwen Guo<sup>1</sup>, Lewei Zeng<sup>1</sup>, Ting Chen<sup>1</sup>, Bohan Yang<sup>1</sup>,  
6 Shupeixiao<sup>1</sup>, Qiongqiong Wang<sup>2</sup>, Zhiyuan Li<sup>3</sup>, Yan You<sup>4</sup>, Shaojun Zhang<sup>5,6,7,8</sup>, and Ye  
7 Wu<sup>5,6,7,8</sup>

8 <sup>1</sup>College of Chemistry and Environmental Engineering, Shenzhen University, Shenzhen 518060, China

9 <sup>2</sup>Department of Atmospheric Science, School of Environmental Studies, China University of  
10 Geosciences, Wuhan 430074, China

11 <sup>3</sup>School of Public Health (Shenzhen), Sun Yat-sen University, Guangzhou 510275, China

12 <sup>4</sup>National Observation and Research Station of Coastal Ecological Environments in Macao, Macao  
13 Environmental Research Institute, Macau University of Science and Technology, Macao SAR 999078,  
14 China

15 <sup>5</sup>School of Environment, State Key Joint Laboratory of Environment Simulation and Pollution Control,  
16 Tsinghua University, Beijing 100084, China

17 <sup>6</sup>State Environmental Protection Key Laboratory of Sources and Control of Air Pollution Complex,  
18 Beijing 100084, China

19 <sup>7</sup>Beijing Laboratory of Environmental Frontier Technologies, School of Environment, Tsinghua  
20 University, Beijing 100084, China

21 <sup>8</sup>Laboratory of Transport Pollution Control and Monitoring Technology, Transport Planning and  
22 Research Institute, Ministry of Transport, Beijing 100028, China

23 *Correspondence to:* Xuan Zheng (x-zheng11@szu.edu.cn)

24 **Abstract:** The advancement of analytical techniques, such as comprehensive two-dimensional gas  
25 chromatography coupled with mass spectrometry (GC×GC-MS), enables the efficient separation of  
26 complex organics. Developing innovative methods for data processing and analysis is crucial to unlock  
27 the full potential of GC×GC-MS in understanding intricate chemical mixtures. In this study, we proposed  
28 an innovative method for the semi-automated identification and quantification of complex organic  
29 mixtures using GC×GC-MS. The method was formulated based on self-constructed mass spectrum  
30 patterns and the traversal algorithms and was applied to organic vapor and aerosol samples collected  
31 from the tailpipe emissions of heavy-duty diesel vehicles and the ambient atmosphere. Thousands of  
32 compounds were filtered, speciated, and clustered into 26 categories, including aliphatic and cyclic  
33 hydrocarbons, aromatic hydrocarbons, aliphatic oxygenated species, phenols and alkyl-phenols, and  
34 heteroatom containing species. The identified species accounted for over 80% of all the eluted  
35 chromatographic peaks at the molecular level. A comprehensive analysis of quantification uncertainty  
36 was undertaken. Using representative compounds, quantification uncertainties were found to be less than  
37 37.67%, 22.54%, and 12.74% for alkanes, polycyclic aromatic hydrocarbons (PAHs), and alkyl-  
38 substituted benzenes, respectively, across the GC×GC space, excluding the first and the last time  
39 intervals. From source apportionment perspective, adamantane was clearly isolated as a potential tracer  
40 for heavy-duty diesel vehicles (HDDVs) emission. The systematic distribution of nitrogen-containing  
41 compounds in oxidized and reduced valences was discussed and many of them served as critical tracers  
42 for secondary nitrate formation processes. The results highlighted the benefits of developing self-  
43 constructed models for the enhanced peak identification, automated cluster analysis, robust uncertainty  
44 estimation, and source apportionment and achieving the full potential of GC×GC-MS in atmospheric  
45 chemistry.

## 46 **1 Introduction**

47 Improved sampling strategies, coupled with innovative measurement techniques, are imperative to  
48 capture the dynamic nature of atmospheric chemistry, particularly in the context of climate change and  
49 health risks (Franklin et al. 2023, Franklin et al. 2022, Huo et al. 2021, Phillips et al. 2018).

50 Comprehensive two-dimensional gas chromatography coupled with mass spectrometry (GC×GC-MS)  
51 has emerged as a powerful tool for compound detection and identification, benefiting from the  
52 combination of two columns with orthogonal selectivity (Alam et al. 2013, Franklin et al. 2022).

53 Despite its capabilities, GC×GC-MS encounters formidable challenges in data analysis, which can be  
54 extremely complicated and demanding. Efforts have been made to handle the deluge of data generated  
55 by GC×GC-MS. Traditionally, mass spectra were deconvoluted and compared to spectra from the  
56 National Institute of Standards and Technology (NIST) library for peak identification with pre-defined  
57 criteria (Guo et al. 2016, Piotrowski et al. 2018). Retention indices (RI) were further introduced to  
58 distinguish homologous compounds with resembling mass spectra. A pioneering and instructive work  
59 for searching criteria to classify GC×GC peaks was published in 2003 (Welthagen, Schnelle-Kreis and  
60 Zimmermann 2003). Welthagen (2003) incorporated the mass fragmentation patterns to classify  
61 compounds in atmospheric aerosol samples. Compounds belonging to the same chemical group related  
62 to one another in the GC×GC space and distributed in a structured pattern. They successfully identified  
63 seven groups of compounds, including alkanes, alkenes and cycloalkanes, alkyl substituted benzenes,  
64 alkyl substituted polar benzenes, hydrated naphthalenes and alkenyl benzenes, alkylated naphthalenes,  
65 and alkane acids, occupying more than 60% of the total peak area. This work set a good example of how  
66 user-defined rules could facilitate the identification of specific compound groups.

67 Recent advances in chemometric tools for GC×GC-MS analysis involving machine learning and deep  
68 learning renovate multi-dimensional chromatography fields (Stefanuto, Smolinska and Focant 2021).  
69 Bendik (2021) developed a programming suite for high-confidence and fast compound identification  
70 using GC×GC coupled with time-of-flight mass spectrometry (TOF-MS) (Bendik et al. 2021). He (2022)  
71 extracted featured mass spectrometric information of the intermediate-volatility and semi-volatile  
72 organic compounds (I/SVOCs) by integrating algorithmic approaches into GC×GC-MS data (He et al.  
73 2022a, He et al. 2022b). A novel pixel-based multiway principal component analysis method was  
74 employed in Song (2023) to identify key tracers during incense burning (Song et al. 2023). Nevertheless,

75 interpreting GC×GC-MS data requires advanced computational tools and expertise, and the investigation  
76 of unknown compounds remains challenging due to the inadequate validation procedures, overreliance  
77 on manual data processing, limited access to computational resources, and the insufficient expertise in  
78 handling complex chromatographic data effectively.

79 Bridging this gap requires further development of sophisticated algorithms and analytical approaches to  
80 unlock the full potential of GC×GC. This study proposes a bottom-up method for cluster analysis and  
81 quantification of organic vapors and aerosols within complex atmospheric mixtures. The scripts were  
82 initiated with the recognition of the common mass spectral features of specific species and were tailored  
83 to a wide range of compound clusters. The scripts were then trained, iterated, and optimized using real  
84 sample data until robust outputs were achieved. The new strategy reduces the ambiguity often associated  
85 with identifying compounds in complex mixtures.

86 The proliferation of heavy-duty diesel vehicles (HDDVs) has raised significant concerns due to their  
87 increasing role in freight transport and in various industrial operations (Yan et al. 2022, Cheng et al.  
88 2022). Despite their low retention rate, HDDVs release substantial amounts of particulate matter,  
89 nitrogen oxides, ammonia, and carbon monoxide into the atmosphere compared with other vehicle types  
90 (Wang et al. 2023, Silva et al. 2023, Chang et al. 2022, Stanimirova et al. 2023, Hamilton and Harley  
91 2021, 2021, Krueve et al. 2014). To address this, gas and aerosol samples were collected from  
92 representative HDDV tailpipes and the ambient environment, then analyzed using GC×GC-MS. The  
93 proposed bottom-up method was employed for a comprehensive analysis of the complex organic  
94 mixtures, resulting in the identification of 26 compound categories, including hydrocarbons in multiple  
95 forms, oxygenated components, and species containing heteroatoms. Over 80% of all the  
96 chromatographic peaks were identified and assigned to a compound cluster using the proposed method,  
97 leaving a minor portion of organic matrix unresolved. Different compound clusters occupied separate  
98 positions in the GC×GC space, and distinctive distribution patterns within diverse samples and their  
99 contribution fractions were revealed. Quantification uncertainties were addressed thoroughly and the  
100 significant potential deviation when using n-alkanes as semi-quantification surrogates was highlighted.  
101 Overall, integrating automated algorithms with GC×GC data analysis holds significant implications for  
102 advancing our understanding of atmospheric chemistry, improving secondary organic aerosol (SOA)  
103 estimation, and guiding the formulation of environmental policies.

## 104 **2 Materials and methods**

### 105 **2.1 Sample collection, treatment, and instrumental analysis**

106 For the collection of HDDVs tailpipe emissions, chassis dynamometer experiments were conducted at  
107 the China Automotive Technology & Research Center (CATARC) in Guangzhou, China. Exhaust  
108 emissions from HDDVs were diluted in a constant volume sampler (CVS, CVS-ONE-MV-HE, Horiba),  
109 following the China heavy-duty commercial vehicle test cycle for tractor trailers (CHTC-TT) driving  
110 cycles. Two HDDVs equipped with the selective catalytic reduction (SCR) system were recruited. The  
111 two HDDVs met the China IV national emission standard and were manufactured in 2021. More  
112 information is summarized in Table S1. The average temperature in the sampling train was precisely  
113 controlled at 47 °C, while airflow, relative humidity, and pressure were monitored simultaneously. The  
114 speed trace and characteristics of CHTC-TT are shown in Figure S1.

115 Gaseous exhausts were collected using two adsorbent thermal desorption (TD) tubes in series (Tenax  
116 TA, C1-AXXX-5003, Markes International) after passing through a Teflon filter. Particulate exhausts  
117 were deposited on a 47 mm quartz filter (Grade QM-A, Whatman). Ambient PM<sub>2.5</sub> filter samples were  
118 collected on the rooftop of a 5-story building on the campus of Shenzhen University (22.60°N, 114.00°E)  
119 during November 2023 in western Shenzhen, approximately 25 m above the ground. The sampling site  
120 was surrounded by campus, residential areas, greenbelts, and a golf park, as shown in Figure S2. Previous  
121 studies demonstrated that the PM<sub>2.5</sub> concentration in this area represented the average pollution scheme  
122 in Shenzhen (Huang et al. 2018, Yu et al. 2020). The sampling strategy followed a regular schedule of  
123 one 24-h sample every day using a high-volume sampler (Th-1000c II, Wuhan Tianhong Environmental  
124 Protection Industry Co., Ltd). In total, 55 TA tube samples (including 11 field blank samples), 20 HDDV  
125 aerosol samples (including 3 field blank samples), and 6 ambient aerosol samples (including one blank  
126 sample) were collected. The list of ambient samples and the relevant PM concentrations are listed in  
127 Table S2. The sorbent tubes were well sealed and stored dry at room temperature, and quartz filters were  
128 frozen at -18 °C before analysis. All sampling materials were pre-baked thoroughly to remove potential  
129 carbonaceous contamination.

130 TD samples were injected with 2 μL of deuterated internal standard (IS) mixing solution through a mild  
131 N<sub>2</sub> blow (CSLR, Markes International). The list of deuterated IS is shown in Table S3. A precise portion  
132 of 1 cm<sup>2</sup> (1 cm × 1 cm) filter sample was isolated and cut into strips. They were spiked with 2 μL of IS

133 mixing solution and inserted into a passivated quartz tube. All TD samples and quartz tubes were loaded  
134 onto a thermal desorption autosampler (ULTRA-xr, Markes International), thermally desorbed (UNITY-  
135 xr, Markes International), and subjected to GC×GC separation (Agilent 8890, Agilent Technologies;  
136 Solid State Modulator1810, J&X Technologies) and mass spectrometry detection (Agilent 5977B,  
137 Agilent Technologies).

138 The thermal desorption system heated the TD tubes to 320 °C (quartz tubes to 330 °C) for 20 min, while  
139 the trap remained at 20 °C. Following tube desorption, the trap temperature was raised to 330 °C (340 °C  
140 for quartz tubes) for 5 min at the maximum heating rate, and the vaporized analytes were purged into the  
141 1<sup>st</sup> GC column with a desorb split flow of 6 mL/min. Separation of the analytes was carried out using a  
142 DB-5ms capillary column (30 m × 0.25 mm × 0.25 μm, Agilent Technologies) as the primary column  
143 and a DB-17ms capillary column (1.2 m × 0.18 mm × 0.18 μm, Agilent Technologies) as the secondary  
144 column. The modulation column consisted of a VF-1ms capillary column (0.7 m × 0.25 mm × 0.10 μm,  
145 Agilent Technologies) connecting to the 1<sup>st</sup> column and an Ultimate Plus deactivated fused silica tubing  
146 (0.6 m × 0.25 mm, Agilent Technologies) connecting to the 2<sup>nd</sup> column.

147 Initially, the GC oven was set at 50 °C for 3 min, followed by a gradual increase at a rate of 5 °C/min  
148 until it reached 310 °C, where it was maintained for an additional 5 min. The entry and exit hot zones  
149 were set +10 °C higher than the GC oven temperature, while the trap zone was maintained at -50 °C. The  
150 modulation cycle had a period of 4 s. Carrier gas flow was set at 1.2 mL/min. The MS had an integer  
151 resolution and was conducted in electron impact positive (EI+) mode (70 eV). It was operated over a  
152 range of 20–350 amu, and the temperature of the transfer line, ion source, and MS quadrupole was 300 °C,  
153 250 °C, and 170 °C, respectively.

## 154 **2.2 Data collection, alignment, and parsing**

155 GC×GC-MS data acquisition was performed using Enhanced MassHunter (version 10.0, Agilent  
156 Technologies) and SSCenter (version 2.4.0.0, J&X Technologies). All data utilized to develop and test  
157 the scripts were processed by Canvas Browser (version 2.5, J&X Technologies), which included baseline  
158 correction, mass spectra deconvolution, and peak smoothing. Baseline correction and peak smoothing  
159 enhanced the signal-to-noise ratio (S/N) and improved overall data quality.

160 Chromatographic peaks were filtered using the following criteria: baseline noise = 150, S/N > 50. For  
161 each individual sample, after isolating all compounds of interest, a peak table was exported with 1<sup>st</sup>

162 retention time (RT) and 2<sup>nd</sup> RT, peak area, peak height, peak width, and deconvoluted mass spectra,  
163 arranged in 1<sup>st</sup> RT sequential order. These quantitative variables were further processed for targeted and  
164 non-targeted “omics”-oriented analysis.

165 As expected, the chromatographic variables experienced RT shifts due to column degradation, routine  
166 maintenance (e.g., cutting column), and system fluctuations (e.g., variations in carrier gas pressure). The  
167 initial tolerance for RT shifts in adaptive cluster matching was set at 1 period of modulation in the 1<sup>st</sup>  
168 dimension and 0.1 s in the 2<sup>nd</sup> dimension. Additionally, a 2D shift cluster consisting of C<sub>16</sub>D<sub>34</sub>, C<sub>24</sub>D<sub>50</sub>,  
169 and C<sub>32</sub>D<sub>66</sub>, was configured, with the merit of correcting 2<sup>nd</sup> RT shift. Data correction or data alignment  
170 is critical for accurate and consistent peak integration.

### 171 **2.3 Algorithmic development**

172 EI spectra are typically characterized by a molecular ion (M<sup>+</sup>) peak plus a collection of fragment ion  
173 peaks. The M<sup>+</sup> may dominate the mass spectrum in some cases (e.g., unsubstituted polycyclic aromatic  
174 hydrocarbons (PAHs)), but more frequently presents at a relatively low intensity. The EI spectra are  
175 highly comparable among different instrument systems and experimental conditions, making them an  
176 excellent measure for identifying compounds. The characteristic ions and their relative intensities depend  
177 on the intrinsic nature of the targeted compounds, necessitating knowledge of basic rules and common  
178 fragmentation routes to interpret EI mass spectra. Figure 1 illustrates the workflow for establishing  
179 computational strategies for robust and reproducible GC×GC-MS data processing.

180 Functional groups significantly affect the fragmentation patterns observed in mass spectrometry, and  
181 some ions are typical of given structures. Isotopic peaks (e.g., hydrogen and chlorine) provide additional  
182 information about the molecules (Du and Angeletti 2006, Fernandez-de-Cossio et al. 2004). These pieces  
183 of information form the foundation for building up the model for cluster analysis, which is addressed in  
184 greater detail in the supporting information (S1). These indicative reaction schemes have been  
185 incorporated into the model development. Each critical step of model construction and validation is  
186 described thoroughly. The quantitative variables in the data alignment table, combining the  
187 chromatographic and MS information, are properly exploited and determine the overall speciation  
188 capacities.

189 Traditionally, compound identification relies on the electron ionization-based fragmentogram and the  
190 deconvoluted mass spectra. Empirically, one-by-one compound identification can be greatly intervened

191 by neighbouring peaks, especially those with similar structures, and introduce considerable uncertainties.  
192 A good example is the assignment of homologous *n*-alkanes, of which the fragmentograms bear a close  
193 resemblance (Figure S6). In such cases, the similarity score (the measure of similarity between the  
194 observed mass spectrum and the NIST library hit) could be erroneously inflated to 850 (out of 999) or  
195 higher. In contrast, cluster analysis involves the comprehensive analysis of a specific type of compounds  
196 on a large scale, aiming to provide a holistic understanding of the distribution and transformation of the  
197 specific compound cluster being investigated.

198 Due to the complexity and remarkable peak capacities, sophisticated and detailed scripts for cluster  
199 identification were constructed. Heteroatom-containing species, e.g., amides and amines, were carefully  
200 examined. The scripts began by recognition of the common mass spectra features of compound cluster  
201 of interest and are addressed in more details in the following descriptive framework:

- 202 1. The Boolean value of characteristic ions.
- 203 2. The intensity sequence of abundant ions in the whole spectra.
- 204 3. The retention time window restriction for certain compound groups.
- 205 4. The pattern of mass spectrometry variation with the increased number of substituents or the  
206 extension of the carbon chain.
- 207 5. An iteration framework that involved repetitive cycles among all the tested samples.

208 The scripts were then trained, iterated, and optimized incorporating real sample data, and the parameters  
209 were adjusted accordingly until a robust output was achieved. The extractor function built into the Canvas  
210 software was activated, and all the scripts were imported to facilitate automated cluster analysis. The  
211 scripts parsed all the files in the given directory into the required structure and generated three reports in  
212 the form of .pdf, .csv, and .bmp. The .csv file contained key information including the compound name,  
213 compound cluster, 1<sup>st</sup> and 2<sup>nd</sup> RTs, and peak area (based on total ion current (TIC)).

```
For (i = 1 to m) # m equals the number of all tested samples.  
  Load the sample  
  Peak identification  
  Baseline correction  
  Mass spectra deconvolution  
  Peak smoothing  
  For (j = 1 to 26) # In total, 26 compound clusters were isolated with high accuracy and  
  repeatability.  
    Execute the extraction rule of cluster (j)
```

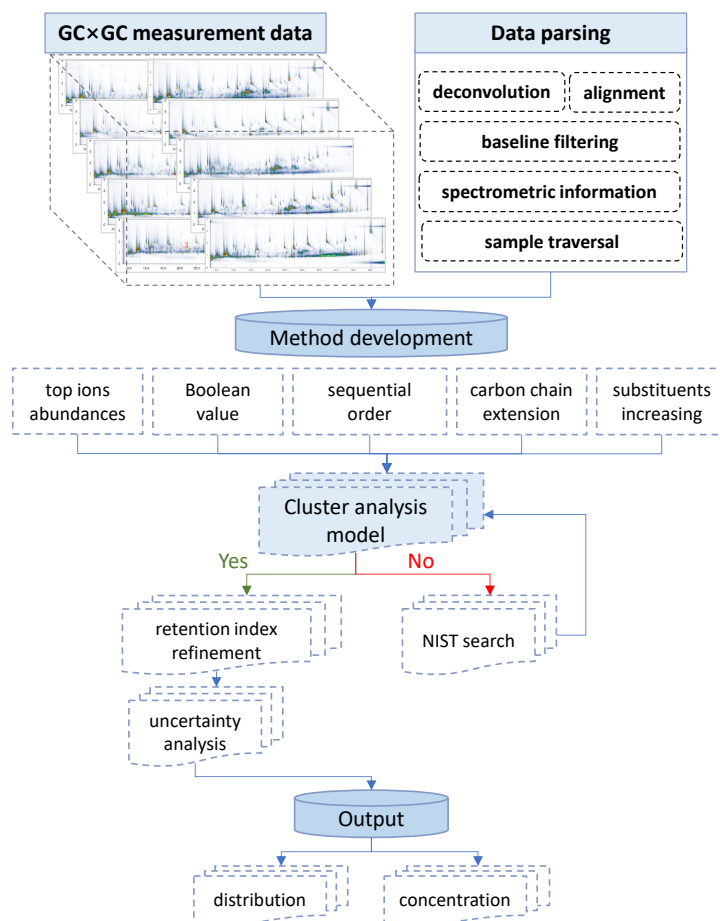


```
Export peak number, 1st RT and 2nd RT, peak area, peak height, peak width, and deconvoluted
mass spectra
Next j
Next i
```

214 Once exported, the peaks were further processed for quantification/semi-quantification following the  
215 steps below. First, calibration curves were prepared by spiking different volumes of the standard solution  
216 mixture onto the blank TD tubes and blank filters, respectively. Peak area ratios, i.e., peak area of  
217 authentic standards over that of the internal standards, were used to build the linear relationship, with the  
218 merit of correcting system fluctuations. The selection of authentic standards prioritized their wide  
219 distribution across the entire chromatogram space, ranging from high to low volatility and weak to strong  
220 polarity, and meanwhile encompassing a broader range of functional groups and heteroatoms. The  
221 distribution and performance of all authentic standards are summarized in Table S4 and Figure S7.  
222 Second, for the un-quantified peaks, their compiled information (X, Y, Z) corresponding to (1<sup>st</sup> RT, 2<sup>nd</sup>  
223 RT, compound cluster) is looped through the list of all authentic standards in the following descriptive  
224 pseudo-codes until the optimal authentic standard to semi-quantify the target peak is exported. It should  
225 be emphasized that the un-quantified peak and the corresponding authentic standard to semi-quantify it  
226 must belong to the same group due to their physicochemical similarities.

```
For (i = 1 to n) # n equals the number of authentic standards and is a known variable.
  If (ZM = Zi) # M is the un-quantified peak and i refers to the authentic standard that is selected in a
  certain loop.
    Ai = Min (an array of ((XM - Xi)2 + (YM - Yi)2)) # This sentence dose not conform to the
    grammar rule of Visual Basic for Applications in Excel, and it is for illustrative purposes only.
    Export Zi, (Xi, Yi, Zi), its peak area, and its linear calibration relationship.
  End if
Next i
```

227



228

229

**Figure 1. Flow diagram illustrating the multistep data processing for establishing computational strategies for cluster analysis and quantification of organic vapors and aerosols using GCxGC-MS data.**

230

#### 231 2.4 Quality assurance/control and uncertainty evaluation

232

It is common for thermal decomposition to occur in analytical methods involving heating processes, potentially leading to the erroneous detection of compounds that are either not present in real samples or present in low concentrations. Such artifacts need careful scrutiny, and the availability of authentic standards covering the GCxGC space range is essential for validation. Nevertheless, the possibility that some observed analytes are decomposition products cannot be entirely ruled out. Peaks of ISs were traced across all samples to monitor the variations across several modules, and the results are presented in Figure S8. Excellent stability was clearly observed, demonstrating the robustness of the testing system. Strong linear correlations were achieved for this set of authentic standards between the peak area ratio and the spiked mass, with Pearson's R ranging from 0.97 to 0.99.

240

241 **3 Results and discussion**

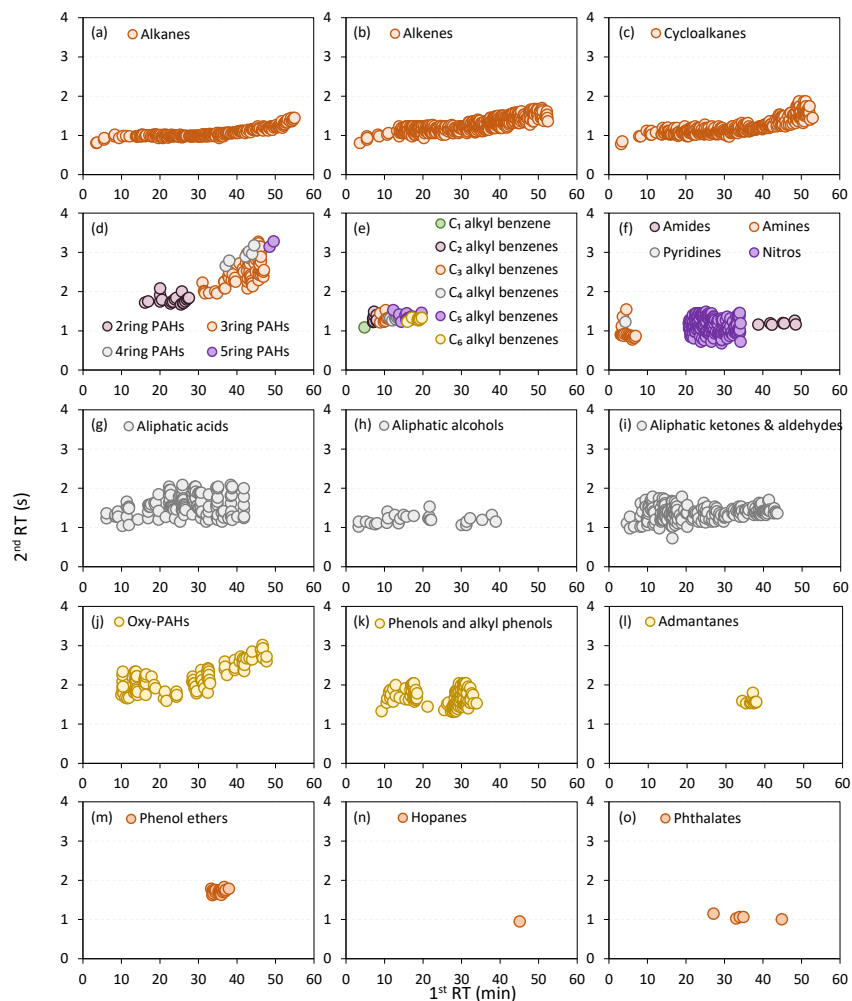
242 **3.1 Overall performance of the algorithm and compound identification**

243 The optimization of component identification remains challenging, and this work involves converting  
244 known chemical compounds into molecular descriptors and utilizing cluster analysis to predict the  
245 relationship between these descriptors and structural information. After continuous trials to improve  
246 reliability and data processing speed, a final solution of 26 compound clusters stands out with high  
247 accuracy and repeatability:

- 248 – Aliphatic hydrocarbons, including *n*-/*i*-alkanes and alkenes
- 249 – Cycloalkanes
- 250 – Alkyl-substituted benzenes, including C<sub>1</sub> – C<sub>6</sub> alkyl-substituted benzenes
- 251 – Adamantanes
- 252 – Hopanes
- 253 – 2 – 5 ring PAHs
- 254 – Acids
- 255 – Aliphatic alcohols
- 256 – Aliphatic aldehydes and ketones
- 257 – Oxy-PAHs
- 258 – Phthalates
- 259 – Phenols and alkyl-substituted phenols
- 260 – Phenol ethers
- 261 – Amides
- 262 – Amines
- 263 – Pyridines
- 264 – Nitro compounds, including organic nitrates and organic nitrites

265 Validation of the model output using field diesel samples has been conducted and has shown high  
266 estimation accuracy and integrity. Generally, over 82% of the peaks have been successfully classified  
267 and assigned to the corresponding compound groups, and their distribution in an example GC×GC plot  
268 is shown in Figure 2. To confirm the tentatively identified heteroatom groups, their raw chromatogram,  
269 mass spectra, and chemical structures of representative species are displayed in Figures S10-S16. Less

270 than 18% of the chromatographic peaks were identified as unresolved components. Aliphatic  
271 hydrocarbons were generally located in the lowest positions in the GC×GC space, except for column  
272 bleedings (Figure 2a-c and Figure S9), and their 2<sup>nd</sup> RT drifted less than 1s from the far-left to the far-  
273 right side. Nitrogen-containing compounds in oxidized and reduced valences, including amides, amines,  
274 pyridines, and nitro compounds, were resolved simultaneously under respective filtering rules and  
275 occupied slightly higher positions in the GC×GC space (Figure 2f). Amines and pyridines, being more  
276 volatile, eluted at early stages, whereas nitro compounds and amides eluted at middle and late stages  
277 sequentially. Due to their high volatility, C<sub>2</sub>-C<sub>6</sub> alkyl-substituted benzenes also appeared at the beginning  
278 of the GC×GC space and predominantly partitioned into the gas phase. Their 2<sup>nd</sup> RTs were comparable  
279 to those of pyridines and amides, with negligible drift in 2<sup>nd</sup> RT. Aliphatic oxygen-containing compounds,  
280 including acids, alcohols, and ketones, were found to be in the middle region and covered a wide  
281 volatility range. These aliphatic oxygen-containing compounds affect the acidity of the atmosphere,  
282 participate in aqueous phase reactions, and contribute significantly to the formation of SOA (Cope et al.  
283 2021, Xu et al. 2022). Phenols with one or more hydroxyl groups attached to an aromatic benzene ring  
284 were observed in the middle of the GC×GC space. Oxy-PAHs and PAHs were present in the upper-  
285 middle of the GC×GC space, with their volatility range extending towards the low volatility end. A clear  
286 trend tilting towards the upper right corner was observed, suggesting that aromaticity plays a significant  
287 role in the retention in the secondary dimension.



288

289 **Figure 2.** The distribution of the 26 compound groups in an example GC×GC plot. For clear visualization,  
 290 different compound groups are displayed separately, except for 2–5 ring PAHs, C<sub>2</sub>–C<sub>6</sub> alkyl-substituted  
 291 benzenes, and nitrogen-containing species. Nitro compounds include organic nitrates and organic nitrites,  
 292 due to the co-existence of the characteristic ions at *m/z* 30 (NO<sup>+</sup>) and *m/z* 46 (NO<sub>2</sub><sup>+</sup>).

### 293 3.2 Estimation of the uncertainty associated with the (semi-) quantification

294 We conducted a systematic evaluation of the model output, and the results are shown in Figure 3 and  
 295 Figure 4. To address this issue comprehensively and accurately, we selected three types of standards  
 296 including C<sub>7</sub>–C<sub>37</sub> *n*-alkanes, C<sub>2</sub>–C<sub>6</sub> alkyl-substituted benzenes, and 2–4 ring PAHs, representing a full  
 297 range of polarities and functionalities. The quantification deviation was computed according to the  
 298 principles of the model. Chromatographic peaks were quantified either by their authentic standards or  
 299 the surrogates within the same compound category after being classified into one of the 26 compound  
 300 classes. For example, if the mass spectrum of a chromatographic peak resembled the pattern of the  
 301 compound class of alkanes, it would be assigned to the alkane group and quantified by its authentic  
 302 standard if available, or by the *n*-alkane (*n*-alkane serving as the semi-quantification surrogate in this

303 case) that was closest to it spatially. Similarly, if the mass spectrum of a chromatographic peak followed  
304 the pattern of C<sub>x</sub> alkyl-substituted benzenes, it would be assigned to the C<sub>x</sub> alkyl-substituted benzene  
305 group and quantified by its authentic standard if available, or by the alkyl-substituted benzene (with  
306 alkyl-substituted benzenes serving as the semi-quantification surrogate) that was closest to it spatially.  
307 In light of this explanation, the deviation of the slopes of the calibration curves of any pair of adjacent  
308 authentic standards within the same compound category was computed to represent the ceiling of the  
309 semi-quantification uncertainty. Uncertainties are calculated using the following Eq. (1):

$$310 \text{ Uncertainty (\%)} = \frac{\text{Abs}(S_p - S_s)}{\text{Smaller}(S_p, S_s)} * 100 \quad (1)$$

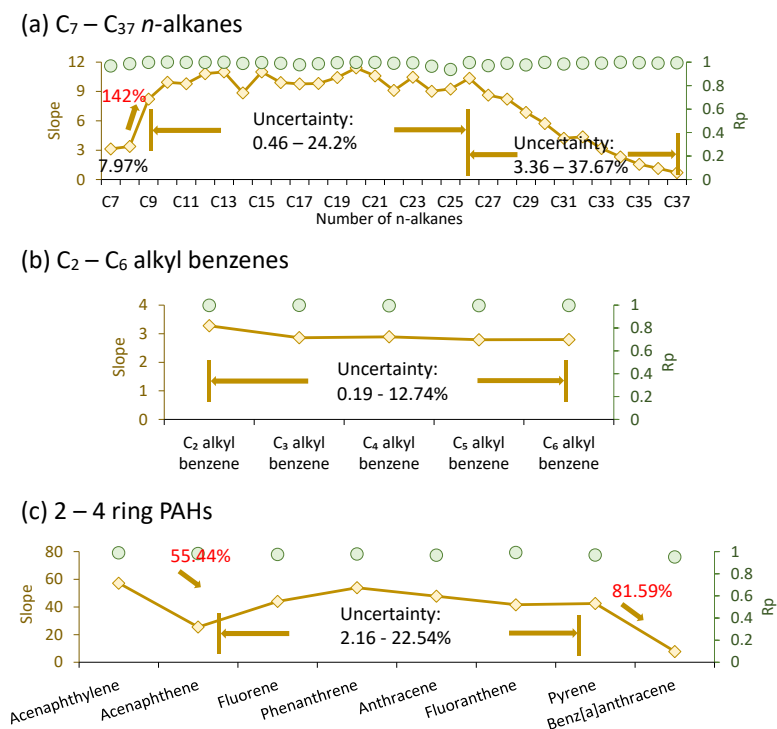
311 where  $S_p$  and  $S_s$  are the slopes of the previous and subsequent compounds, respectively.

312 The slopes increased rapidly from 3.13 (C<sub>7</sub> *n*-alkane) to 8.21 (C<sub>9</sub> *n*-alkane), fluctuated slightly from 8.85  
313 to 11.8 in the range of C<sub>9</sub> to C<sub>27</sub> *n*-alkanes, and decreased gradually after C<sub>28</sub> *n*-alkane to the end of C<sub>37</sub>  
314 *n*-alkane. Throughout the volatility range of C<sub>9</sub>–C<sub>37</sub> *n*-alkanes, uncertainties were less than 37.67%,  
315 except for one interval between C<sub>8</sub> and C<sub>9</sub> *n*-alkanes, where the quantification deviation reached 142%.  
316 A similar trend was observed for PAHs, with uncertainties less than 22.54%, except for the first and last  
317 intervals, where the quantification deviations were 55.44% and 81.59%, respectively, as shown in Figure  
318 3. Stable responses of C<sub>2</sub>–C<sub>6</sub> alkyl-substituted benzenes were monitored, and the uncertainties were less  
319 than 12.74%. In other words, for any given peak, it would be quantified or semi-quantified by one  
320 authentic standard, and the upper limit of quantification uncertainty, originating from any pair of adjacent  
321 authentic standards, was as discussed earlier.

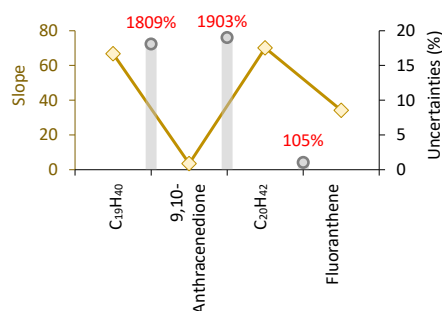
322 It is reasonable that the uncertainty ranges of alkyl-substituted benzenes were less than those of *n*-alkanes  
323 and PAHs, given that alkyl-substituted benzenes eluted early in the front half, whereas alkanes and PAHs  
324 covered the entire volatility range. These trends illustrated that the responses of GC×GC to the analysts  
325 were sensitive to the volatility distribution, with accurate quantification being more reliable in the middle  
326 region. This also highlighted the utility of introducing more authentic standards and the benefits of  
327 enriching compound categories. It can be speculated that the quantification uncertainty would be further  
328 reduced with the addition of more standard compounds.

329 Furthermore, we explored the uncertainty estimation of dividing the whole chromatogram into bins based  
330 on retention time, and all the species in the same bin were quantified, referring to the mass-to-signal  
331 responses of the C<sub>n</sub> *n*-alkanes (Zhao et al. 2015, Zhao et al. 2014). This approach corrected the signal

332 variation of hydrocarbons in the GC-MS and was widely adopted for quantifying unresolved complex  
 333 mixtures (UCMs) (Shen et al. 2023, Zhao et al. 2022). We chose four types of standards belonging to  
 334 different compound categories with similar 1<sup>st</sup> RTs and different 2<sup>nd</sup> RTs, including C<sub>19</sub>H<sub>40</sub> (1<sup>st</sup> RT = 34.6  
 335 min, 2<sup>nd</sup> RT = 1.03 s), 9,10-anthracenedione (1<sup>st</sup> RT = 36.07 min, 2<sup>nd</sup> RT = 3.85 s), C<sub>19</sub>H<sub>40</sub> (1<sup>st</sup> RT = 36.54  
 336 min, 2<sup>nd</sup> RT = 1.07 s), and fluoranthene (1<sup>st</sup> RT = 37.00 min, 2<sup>nd</sup> RT = 3.04 s), and assessed the deviation  
 337 of slopes between each pair of the standards. Results shown in Figure 4 indicate that the deviation  
 338 between the three pairs of standards was 1809% (C<sub>19</sub> *n*-alkane vs. 9,10-anthracenedione), 1903% (9,10-  
 339 anthracenedione vs. C<sub>20</sub> *n*-alkane), and 105% (C<sub>20</sub> *n*-alkane vs. fluoranthene), respectively. Quantitative  
 340 errors in measuring unidentified chromatographic peaks using *n*-alkanes responses could reach three  
 341 orders of magnitude, especially for oxygen-containing species. Errors in quantifying aromatic  
 342 components, e.g., PAHs, also exceeded 100% in some cases.



343  
 344 **Figure 3. Slope and Pearson correlation variation of (a) C<sub>7</sub>–C<sub>37</sub> *n*-alkanes, (b) C<sub>2</sub>–C<sub>6</sub> alkyl-substituted**  
 345 **benzenes, and (c) 2–4 ring PAHs. Brown diamond dots represent slopes of different species and are referenced**  
 346 **to the left axis. Green circles denote the Pearson correlation of individual species and are referenced to the**  
 347 **right axis. Pearson correlation values for *n*-alkanes, C<sub>2</sub>–C<sub>6</sub> alkyl-substituted benzenes, and 2–4 ring PAHs**  
 348 **range from 0.936 to 0.999, 0.994 to 0.998, and 0.952 to 0.992, respectively. Uncertainties are computed using**  
 349 **the equation provided in the main text.**



350

351 **Figure 4. Slopes and uncertainty estimation for example compounds with close 1<sup>st</sup> RTs and different 2<sup>nd</sup> RTs:**  
 352 **C<sub>19</sub>H<sub>40</sub> (1<sup>st</sup> RT = 34.6 min, 2<sup>nd</sup> RT = 1.03 s), 9,10-anthracenedione (1<sup>st</sup> RT = 36.07 min, 2<sup>nd</sup> RT = 3.85 s), C<sub>20</sub>H<sub>42</sub>**  
 353 **(1<sup>st</sup> RT = 36.54 min, 2<sup>nd</sup> RT = 1.07 s), and fluoranthene (1<sup>st</sup> RT = 37.00 min, 2<sup>nd</sup> RT = 3.04 s). Brown diamond**  
 354 **dots represent the slopes of different species and are referenced to the left axis. Gray bars denote the**  
 355 **uncertainty estimation for these compounds and are referenced to the right axis.**

### 356 3.3 Cluster analysis in organic vapor and aerosol samples

357 The model was applied to organic vapor samples from HDDV tailpipe emissions (referred to as HDDV  
 358 vapors), aerosol samples from HDDV tailpipe emissions (referred to as HDDV aerosols), and  
 359 atmospheric aerosol samples (referred to as ambient aerosols) for cluster analysis. The results are shown  
 360 in Figure S17, which displays the distribution of the top few species with a contribution fraction  
 361 exceeding 5%, and in Figure 5, which shows the mass stacking. Overall, the identified chromatographic  
 362 peaks accounted for 85%, 82%, and 99% for HDDV vapors, HDDV aerosols, and ambient aerosol  
 363 samples, respectively. The unidentified peaks were less than 20% and are addressed in greater detail in  
 364 the supporting information (S2).

365 Distinct cluster distribution features can be extracted. For ambient aerosol samples, aliphatic ketones  
 366 were the most abundant cluster, contributing to 27% of all the peak signals, followed by alkanes and  
 367 alkenes. A notable fraction of 15.2% of nitro compounds was observed exclusively in ambient samples,  
 368 indicating significant secondary nitrate formation processes under atmospheric conditions. Aliphatic  
 369 acids and oxy-PAHs were also detected at high levels, with the top six groups accounting for over 95%  
 370 of the total classified peak signals. Minor but non-negligible fractions included cycloalkane, aliphatic  
 371 alcohols, and phenols and alkyl-substituted phenols.

372 Similarly, aliphatic ketones ranked first for HDDV aerosol samples, with mass intensity reaching 46%  
 373 of the total signals, followed by alkanes. Aliphatic alcohols and oxy-PAHs were also detected at high  
 374 levels, and the top four groups accounted for over 88% of the total classified peak signals. Cycloalkanes,



375 amides, phenols and alkyl-substituted phenols, and alkenes were compound clusters with lower  
376 abundance, ranging from 1% to 4%.

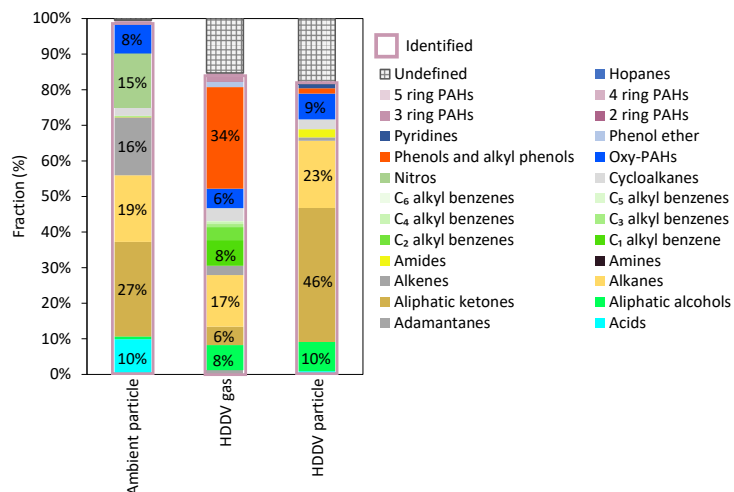
377 For HDDV vapors, the most abundant group was phenols and alkyl-substituted phenols, constituting 34%  
378 of the total peak signals. Compared with previous results, where the most abundant group was reported  
379 to be alkanes (Wang et al. 2022, Alam et al. 2019), the adoption of the innovative model contributed to  
380 resolving the oxygenated fractions and reduced inaccuracies in SOA simulation due to the lack of species  
381 information. The compound cluster is confirmed by 1) the retention time window, including 1<sup>st</sup> RT and  
382 2<sup>nd</sup> RT, and 2) the mass spectra. Detailed information is displayed in Figure S16. The 2<sup>nd</sup> RTs of the  
383 identified phenols and alkyl-substituted phenols range from 1.45 to 1.78 s, well above the hydrocarbon  
384 regions, where the 2<sup>nd</sup> RTs fall within the range of approximately 1.0 to 1.15 s. Their mass spectra also  
385 feature the typical phenol ions at  $m/z = 94, 107, 121, 135, 149, \text{ and } 191$ . Alkanes ranked as the second  
386 top species, followed by C<sub>1</sub> alkyl-substituted benzene. C<sub>1</sub>–C<sub>6</sub> alkyl-substituted benzenes were negligible  
387 in both ambient and HDDV aerosol samples but were present in notable abundance in HDDV vapor  
388 samples. This distribution aligned with their placement in the GC×GC plot, indicating they were  
389 relatively volatile species and partitioned predominantly into the gas phase. Oxy-PAHs and aliphatic  
390 ketones contributed 6% of the total identified peak intensities, followed by minor fractions including C<sub>2</sub>  
391 alkyl-substituted benzene, cycloalkanes, and alkenes.

392 The model output illustrates the overall distribution of compound clusters in various gas and aerosol  
393 samples, providing comparative insights. Carboxylic acids indicated a higher oxidation state than other  
394 compound clusters and were exclusively observed at a notable level in ambient samples compared with  
395 “freshly emitted” source samples. The oxidation state of dominant compounds in HDDV samples was  
396 comparatively low. For example, a significant ketone fraction was observed in HDDV samples, with the  
397 majority partitioning into the aerosol phase due to the long carbon chain skeleton and thus low volatility.  
398 Phenols and alkyl-substituted phenols were the leading species in HDDV gas samples. He (2022)  
399 reported that the oxygenated I/SVOCs accounted for over 20% of the total I/SVOCs mass in HDDV  
400 tailpipe emissions (He et al. 2022a). With the refinement and improvement of model performance, e.g.,  
401 further splitting mixed mixtures, the oxygenated fraction was elevated to over 50%.

402 This study highlighted the systematic presence and distribution of nitrogen-containing compounds in  
403 both oxidized valences (including nitro compounds) and reduced valences (including amides, amines,  
404 and pyridines). Among them, amines and amides were key precursors for new particle formation

405 processes in a polluted atmosphere (Saeki et al. 2022, Cai et al. 2021), and pyridines, with the nitrogen  
406 atom in the aromatic ring, were readily dissolved in water, participating in the global nitrogen cycle in  
407 ecosystems (Kosyakov et al. 2020). Nitro compounds, which include a wide range of organic compounds  
408 with NO or NO<sub>2</sub> substituents, served as critical tracers for secondary nitrate formation processes. Amines  
409 and pyridines were volatile species occupying the early section of the GC×GC space, while nitro  
410 compounds and amides were distributed in the middle and rear space. Individual nitrogen-containing  
411 species were present at trace levels under atmospheric conditions and were difficult to detect. Moreover,  
412 authentic standards or high-resolution mass spectrometry were required to identify and quantify each  
413 compound (Zhang et al. 2018). With the establishment of an algorithmic solution, we were able to  
414 conduct a comprehensive scan of nitrogen-containing compound clusters.

415 In addition to common features, specific compounds were identified in separate samples and could  
416 potentially serve as markers or tracers for primary emissions. Adamantane and its derivatives, with the  
417 fusion of three cyclohexane rings (chemical structure and mass spectrum shown in Figure S18a), are  
418 natural products in petroleum (Stout and Douglas 2004). They were volatile and had previously been  
419 isolated using GC×GC-ToF-MS in crude oil (Wang et al. 2013). Adamantanes were observed in HDDV  
420 vapor samples, contributing 1.4% to the identified peaks. Hopane (chemical structure and mass spectrum  
421 shown in Figure S18b) is also a natural product in petroleum and bitumen and serves as an important  
422 marker for vehicle emissions due to its persistency and stability (He et al. 2022b, Wong et al. 2021).  
423 Hopane was reported to survive heat treatment up to 460 °C and was exclusively detected in HDDV  
424 aerosol samples, with an intensity fraction of 0.3% (Wu and Geng 2016).



425

426 **Figure 5. Fractional distributions of different compound clusters in ambient aerosol samples, HDDV tailpipe**  
 427 **vapors, and HDDV tailpipe aerosols. Numbers labelled on each column represent the fractions of the top few**  
 428 **groups in different samples. Identified clusters are outlined in light purple.**

429 **4 Conclusions and outlook**

430 We presented an innovative method for optimizing the separation and identification of organic vapors  
 431 and aerosols, focusing on establishing molecular descriptors and cluster analysis algorithms. The model  
 432 outputs were validated using field samples with high accuracy and integrity. Less than 20% of the peaks  
 433 were unresolved components. The retention patterns of various compound groups and their distribution  
 434 in the GC×GC plot were resolved, and the influence of functional groups on fragmentation was  
 435 thoroughly addressed. We also provided a comprehensive analysis of the quantification uncertainties of  
 436 this new approach and highlighted the significant quantitative errors when using *n*-alkanes as semi-  
 437 quantification surrogates. This model was applied to various types of field samples, and the results  
 438 revealed distinctive distribution patterns of compound clusters and contribution fractions, providing  
 439 valuable insights into the compositions of organic vapors and aerosols, and offering potential markers  
 440 for specific emission sources.

441 Compound speciation in atmospheric chemistry continues to be a dynamic and challenging field.  
 442 Speciated compounds enable models to consider the diversity of organic species and dynamic chemical  
 443 transformations in the atmosphere, contributing to more accurate SOA simulation results. This approach  
 444 also allows for a more refined description of the dispersion of pollutants, thereby assisting in the  
 445 development of localized air quality management strategies as we strive for a more accurate and  
 446 comprehensive understanding of atmospheric chemistry.

447 **Supplement link:**

448 **Data availability**

449 The measurement data and codes used in this study are available on request.

450 **Author contribution:**

451 X.H.: Conceptualization, formal analysis, model development, data validation, writing—original draft,  
452 funding acquisition; X.Z.: Writing—reviewing and editing, project administration, supervision, funding  
453 acquisition; S.G: Experiment; L.Z., T.C., B.Y., and S.X.: Experiment; Q.W., Z.L., Y.Y., S.Z., and Y.W.:  
454 Data validation, writing—reviewing and editing.

455 **Competing interests:**

456 The authors declare that they have no conflict of interest.

457 **Acknowledgements**

458 The authors acknowledge the financial support of the National Natural Science Foundation of China  
459 (Grant No. 42105100 and 42261160645), Scientific Research Fund at Shenzhen University (868-  
460 000001032089 and 827-000907), Shenzhen Outstanding Science and Technology Innovation Program  
461 (RCBS20231211090534051), and Macao Science and Technology Development Fund (0023/2022/AFJ).

462 **References**

463 2021. Ministry of Ecology and Environment of the People's Republic of China. China Mobile Source  
464 Environmental Management Annual Report 2021.  
465 Alam, M. S., C. E. West, A. G. Scarlett, S. J. Rowland & R. M. Harrison Application of 2D-GCMS  
466 reveals many industrial chemicals in airborne particulate matter. *Atmospheric Environment*, 65, 101-111,  
467 2013.  
468 Alam, M. S., S. Zeraati-Rezaei, H. M. Xu & R. M. Harrison Characterization of Gas and Particulate  
469 Phase Organic Emissions (C<sub>9</sub>-C<sub>37</sub>) from a Diesel Engine and the Effect of Abatement Devices.  
470 *Environmental Science & Technology*, 53, 11345-11352, 2019.  
471 Bendik, J., R. Kalia, J. Sukumaran, W. H. Richardot, E. Hoh & S. T. Kelley Automated high confidence  
472 compound identification of electron ionization mass spectra for nontargeted analysis. *J Chromatogr A*,  
473 1660, 462656, <https://www.ncbi.nlm.nih.gov/pubmed/34798444>, 2021.

474 Cai, R. L., C. Yan, D. R. Worsnop, F. Bianchi, V. M. Kerminen, Y. C. Liu, L. Wang, J. Zheng, M.  
475 Kulmala & J. K. Jiang An indicator for sulfuric acid-amine nucleation in atmospheric environments.  
476 *Aerosol Science and Technology*, 55, 1059-1069, 2021.

477 Chang, Y. H., K. Cheng, Y. Q. Kuang, Q. Y. Hu, Y. Q. Gao, R. J. Huang, C. Huang, W. W. Walters &  
478 M. F. Lehmann Isotopic Variability of Ammonia ( $\delta^{15}\text{N-NH}_3$ ) Slipped from Heavy-Duty Vehicles under  
479 Real-World Conditions. *Environ Sci Technol Letters*, 9, 726-732, 2022.

480 Cheng, S. F., Y. B. Zhao, B. B. Zhang, P. Peng & F. Lu Structural decomposition of heavy-duty diesel  
481 truck emission contribution based on trajectory mining. *Journal of Cleaner Production*, 380, 135172,  
482 2022.

483 Cope, J. D., K. A. Abellar, K. H. Bates, X. Fu & T. B. Nguyen Aqueous Photochemistry of 2-Methyltetrol  
484 and Erythritol as Sources of Formic Acid and Acetic Acid in the Atmosphere. *ACS Earth and Space*  
485 *Chemistry*, 5, 1265-1277, <https://doi.org/10.1021/acsearthspacechem.1c00107>, 2021.

486 Du, P. & R. H. Angeletti Automatic deconvolution of isotope-resolved mass spectra using variable  
487 selection and quantized peptide mass distribution. *Anal Chem*, 78, 3385-3392,  
488 <https://www.ncbi.nlm.nih.gov/pubmed/16689541>, 2006.

489 Fernandez-de-Cossio, J., L. J. Gonzalez, Y. Satomi, L. Betancourt, Y. Ramos, V. Huerta, V. Besada, G.  
490 Padron, N. Minamino & T. Takao Automated interpretation of mass spectra of complex mixtures by  
491 matching of isotope peak distributions. *Rapid Commun Mass Spectrom*, 18, 2465-2472,  
492 <https://www.ncbi.nlm.nih.gov/pubmed/15384131>, 2004.

493 Franklin, E. B., S. Amiri, D. Crocker, C. Morris, K. Mayer, J. S. Sauer, R. J. Weber, C. Lee, F. Malfatti,  
494 C. D. Cappa, T. H. Bertram, K. A. Prather & A. H. Goldstein Anthropogenic and Biogenic Contributions  
495 to the Organic Composition of Coastal Submicron Sea Spray Aerosol. *Environ Sci Technol*, 56, 16633-  
496 16642, <https://www.ncbi.nlm.nih.gov/pubmed/36332100>, 2022.

497 Franklin, E. B., L. D. Yee, R. Wernis, G. Isaacman-VanWertz, N. Kreisberg, R. Weber, H. Zhang, B. B.  
498 Palm, W. Hu, P. Campuzano-Jost, D. A. Day, A. Manzi, P. Artaxo, R. A. F. Souza, J. L. Jimenez, S. T.  
499 Martin & A. H. Goldstein Chemical Signatures of Seasonally Unique Anthropogenic Influences on  
500 Organic Aerosol Composition in the Central Amazon. *Environ Sci Technol*, 57, 6263-6272,  
501 <https://www.ncbi.nlm.nih.gov/pubmed/37011031>, 2023.

502 Guo, Q., J. Yu, K. Yang, X. Wen, H. Zhang, Z. Yu, H. Li, D. Zhang & M. Yang Identification of complex  
503 septic odorants in Huangpu River source water by combining the data from gas chromatography-  
504 olfactometry and comprehensive two-dimensional gas chromatography using retention indices. *Sci Total*  
505 *Environ*, 556, 36-44, <https://www.ncbi.nlm.nih.gov/pubmed/26974564>, 2016.

506 Hamilton, S. D. & R. A. Harley High-Resolution Modeling and Apportionment of Diesel-Related  
507 Contributions to Black Carbon Concentrations. *Environ Sci Technol*, 55, 12250-12260, 2021.

508 He, X., X. Zheng, Y. You, S. Zhang, B. Zhao, X. Wang, G. Huang, T. Chen, Y. Cao, L. He, X. Chang,  
509 S. Wang & Y. Wu Comprehensive chemical characterization of gaseous I/SVOC emissions from heavy-  
510 duty diesel vehicles using two-dimensional gas chromatography time-of-flight mass spectrometry.  
511 *Environ Pollut*, 305, 119284, <https://www.ncbi.nlm.nih.gov/pubmed/35436508>, 2022a.

512 He, X., X. Zheng, S. Zhang, X. Wang, T. Chen, X. Zhang, G. Huang, Y. Cao, L. He, X. Cao, Y. Cheng,  
513 S. Wang & Y. Wu Comprehensive characterization of particulate intermediate-volatility and semi-  
514 volatile organic compounds (I/SVOCs) from heavy-duty diesel vehicles using two-dimensional gas  
515 chromatography time-of-flight mass spectrometry. *Atmos. Chem. Phys.*, 22, 13935-13947,  
516 <https://dx.doi.org/10.5194/acp-22-13935-2022>, 2022b.

517 Huang, X. F., B. B. Zou, L. Y. He, M. Hu, A. S. H. Prévôt & Y. H. Zhang Exploration of PM<sub>2.5</sub> sources  
518 on the regional scale in the Pearl River Delta based on ME-2 modeling. *Atmos. Chem. Phys.*, 18, 11563-  
519 11580, <https://acp.copernicus.org/articles/18/11563/2018/>, 2018.

520 Huo, Y., Z. Guo, Q. Li, D. Wu, X. Ding, A. Liu, D. Huang, G. Qiu, M. Wu, Z. Zhao, H. Sun, W. Song,  
521 X. Li, Y. Chen, T. Wu & J. Chen Chemical Fingerprinting of HULIS in Particulate Matters Emitted from  
522 Residential Coal and Biomass Combustion. *Environ Sci Technol*, 55, 3593-3603,  
523 <https://www.ncbi.nlm.nih.gov/pubmed/33656861>, 2021.

524 Kosyakov, D. S., N. V. Ul'yanovskii, T. B. Latkin, S. A. Pokryshkin, V. R. Berzhonskis, O. V. Polyakova  
525 & A. T. Lebedev Peat burning - An important source of pyridines in the earth atmosphere. *Environ Pollut*,  
526 266, 115109, <https://www.ncbi.nlm.nih.gov/pubmed/32622216>, 2020.

527 Kruve, A., K. Kaupmees, J. Liigand & I. Leito Negative electrospray ionization via deprotonation:  
528 predicting the ionization efficiency. *Anal Chem*, 86, 4822-30,  
529 <https://www.ncbi.nlm.nih.gov/pubmed/24731109>, 2014.

530 Phillips, K. A., A. Yau, K. A. Favela, K. K. Isaacs, A. McEachran, C. Grulke, A. M. Richard, A. J.  
531 Williams, J. R. Sobus, R. S. Thomas & J. F. Wambaugh Suspect Screening Analysis of Chemicals in  
532 Consumer Products. *Environ Sci Technol*, 52, 3125-3135,  
533 <https://www.ncbi.nlm.nih.gov/pubmed/29405058>, 2018.

534 Piotrowski, P. K., B. A. Weggler, D. A. Yoxtheimer, C. N. Kelly, E. Barth-Naftilan, J. E. Saiers & F. L.  
535 Dorman Elucidating Environmental Fingerprinting Mechanisms of Unconventional Gas Development  
536 through Hydrocarbon Analysis. *Anal Chem*, 90, 5466-5473,  
537 <https://www.ncbi.nlm.nih.gov/pubmed/29580048>, 2018.

538 Saeki, K., K. Ikari, H. Yokoi, S. I. Ohira, H. Okochi & K. Toda Biogenic Diamines and Their Amide  
539 Derivatives Are Present in the Forest Atmosphere and May Play a Role in Particle Formation. *ACS Earth  
540 and Space Chemistry*, 6, 421-430, 2022.

541 Shen, X., H. Che, Z. Yao, B. Wu, T. Lv, W. Yu, X. Cao, X. Hao, X. Li, H. Zhang & X. Yao Real-World  
542 Emission Characteristics of Full-Volatility Organics Originating from Nonroad Agricultural Machinery  
543 during Agricultural Activities. *Environ Sci Technol*, 57, 10308-10318,  
544 <https://www.ncbi.nlm.nih.gov/pubmed/37419883>, 2023.

545 Silva, L. F. M., A. R. H. De La Cruz, A. H. M. Nunes & A. Gioda Real-Time Monitoring of Nitrogen  
546 Oxides Emission Factors Using Sensors in the Exhaust Pipes of Heavy Vehicles in the Metropolitan  
547 Region of Rio de Janeiro. *Journal of the Brazilian Chemical Society*, 2023.

548 Song, K., R. Tang, J. Zhang, Z. Wan, Y. Zhang, K. Hu, Y. Gong, D. Lv, S. Lu, Y. Tan, R. Zhang, A. Li,  
549 S. Yan, S. Yan, B. Fan, W. Zhu, C. K. Chan & S. Guo. 2023. Molecular fingerprints and health risks of  
550 home-use incense burning smoke. Copernicus GmbH.

551 Stanimirova, I., D. Q. Rich, A. G. Russell & P. K. Hopke A long-term, dispersion normalized PMF  
552 source apportionment of PM<sub>2.5</sub> in Atlanta from 2005 to 2019. *Atmospheric Environment*, 312, 120027,  
553 2023.

554 Stefanuto, P.-H., A. Smolinska & J.-F. Focant Advanced chemometric and data handling tools for  
555 GC×GC-TOF-MS Application of chemometrics and related advanced data handling in chemical  
556 separations. *TrAC Trends in Analytical Chemistry*, 139, 116251, 2021.

557 Stout, S. A. & G. S. Douglas Diamondoid hydrocarbons - Application in the chemical fingerprinting of  
558 natural gas condensate and gasoline. *Environmental Forensics*, 5, 225-235, 2004.

559 Wang, A. Q., Z. B. Yuan, X. H. Liu, M. L. Wang, J. Yang, Q. E. Sha & J. Y. Zheng Measurement-based  
560 intermediate volatility organic compound emission inventory from on-road vehicle exhaust in China.  
561 *Environmental Pollution*, 310, 2022.

562 Wang, G. L., S. B. Shi, P. R. Wang & T. G. Wang Analysis of diamondoids in crude oils using  
563 comprehensive two-dimensional gas chromatography/time-of-flight mass spectrometry. *Fuel*, 107, 706-  
564 714, 2013.

565 Wang, H., S. J. Zhang, X. M. Wu, Y. F. Wen, Z. H. Li & Y. Wu Emission Measurements on a Large  
566 Sample of Heavy-Duty Diesel Trucks in China by Using Mobile Plume Chasing. *Environ Sci Technol*,  
567 57, 15153-15161, 2023.

568 Welthagen, W., J. Schnelle-Kreis & R. Zimmermann Search criteria and rules for comprehensive two-  
569 dimensional gas chromatography-time-of-flight mass spectrometry analysis of airborne particulate  
570 matter. *J Chromatogr A*, 1019, 233-49, <https://www.ncbi.nlm.nih.gov/pubmed/14650618>, 2003.

571 Wong, Y. K., X. H. H. Huang, Y. Y. Cheng & J. Z. Yu Estimating primary vehicular emission  
572 contributions to PM<sub>2.5</sub> using the Chemical Mass Balance model: Accounting for gas-particle partitioning  
573 of organic aerosols and oxidation degradation of hopanes. *Environmental Pollution*, 291, 118131, 2021.

574 Wu, L. L. & A. S. Geng Differences in the thermal evolution of hopanes and steranes in free and bound  
575 fractions. *Organic Geochemistry*, 101, 38-48, 2016.

576 Xu, B., G. Zhang, O. Gustafsson, K. Kawamura, J. Li, A. Andersson, S. Bikkina, B. Kunwar, A. Pokhrel,  
577 G. Zhong, S. Zhao, J. Li, C. Huang, Z. Cheng, S. Zhu, P. Peng & G. Sheng Large contribution of fossil-  
578 derived components to aqueous secondary organic aerosols in China. *Nat Commun*, 13, 5115,  
579 <https://www.ncbi.nlm.nih.gov/pubmed/36045131>, 2022.

580 Yan, J. Z., G. Wang, S. Y. Chen, H. Zhang, J. Q. Qian & Y. X. Mao Harnessing freight platforms to  
581 promote the penetration of long-haul heavy-duty hydrogen fuel-cell trucks. *Energy*, 254, 124225, 2022.

582 Yu, D., Z. Tan, K. Lu, X. Ma, X. Li, S. Chen, B. Zhu, L. Lin, Y. Li, P. Qiu, X. Yang, Y. Liu, H. Wang,  
583 L. He, X. Huang & Y. Zhang An explicit study of local ozone budget and NO<sub>x</sub>-VOCs sensitivity in  
584 Shenzhen China. *Atmospheric Environment*, 224, 117304,  
585 <https://www.sciencedirect.com/science/article/pii/S1352231020300467>, 2020.

586 Zhang, Y., R. Li, J. Fang, C. Wang & Z. Cai Simultaneous determination of eighteen nitro-polyaromatic  
587 hydrocarbons in PM<sub>2.5</sub> by atmospheric pressure gas chromatography-tandem mass spectrometry.  
588 *Chemosphere*, 198, 303-310, <https://www.ncbi.nlm.nih.gov/pubmed/29421744>, 2018.

589 Zhao, Y., C. J. Hennigan, A. A. May, D. S. Tkacik, J. A. de Gouw, J. B. Gilman, W. C. Kuster, A. Borbon  
590 & A. L. Robinson Intermediate-volatility organic compounds: a large source of secondary organic  
591 aerosol. *Environ Sci Technol*, 48, 13743-137550, <https://www.ncbi.nlm.nih.gov/pubmed/25375804>,  
592 2014.

593 Zhao, Y., N. T. Nguyen, A. A. Presto, C. J. Hennigan, A. A. May & A. L. Robinson Intermediate  
594 Volatility Organic Compound Emissions from On-Road Diesel Vehicles: Chemical Composition,  
595 Emission Factors, and Estimated Secondary Organic Aerosol Production. *Environ Sci Technol*, 49,  
596 11516-11526, <https://www.ncbi.nlm.nih.gov/pubmed/26322746>, 2015.

597 Zhao, Y., D. S. Tkacik, A. A. May, N. M. Donahue & A. L. Robinson Mobile Sources Are Still an  
598 Important Source of Secondary Organic Aerosol and Fine Particulate Matter in the Los Angeles Region.  
599 *Environ Sci Technol*, 56, 15328-15336, <https://www.ncbi.nlm.nih.gov/pubmed/36215417>, 2022.

600 For Table of Contents Only

