

Accepted Manuscript

Structure/reaction directed analysis for LC-MS based untargeted analysis

Miao Yu, Mariola Olkowicz, Janusz Pawliszyn

PII: S0003-2670(18)31304-7

DOI: <https://doi.org/10.1016/j.aca.2018.10.062>

Reference: ACA 236368

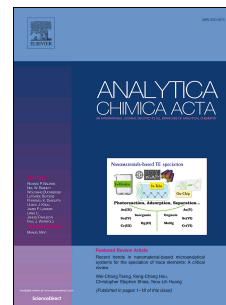
To appear in: *Analytica Chimica Acta*

Received Date: 19 August 2018

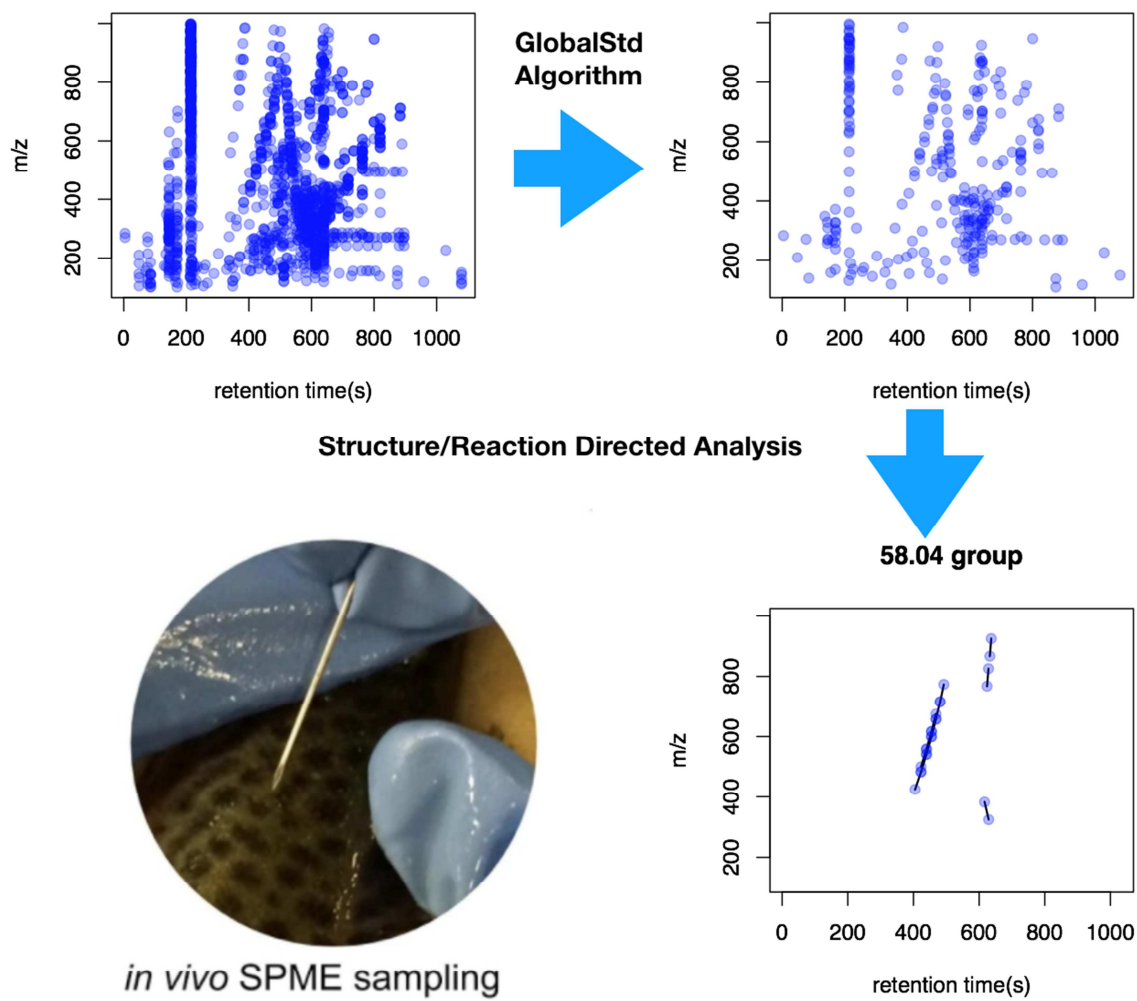
Accepted Date: 25 October 2018

Please cite this article as: M. Yu, M. Olkowicz, J. Pawliszyn, Structure/reaction directed analysis for LC-MS based untargeted analysis, *Analytica Chimica Acta*, <https://doi.org/10.1016/j.aca.2018.10.062>.

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.



The final publication is available at Elsevier via <https://doi.org/10.1016/j.aca.2018.10.062>. © 2018. This manuscript version is made available under the CC-BY-NC-ND 4.0 license <http://creativecommons.org/licenses/by-nc-nd/4.0/>



ACCEPTED

1 **Structure/reaction directed analysis for LC-**
2 **MS based untargeted analysis**

3 Miao Yu¹, Mariola Olkowicz¹ and Janusz Pawliszyn^{1*}

4 ¹Department of Chemistry, University of Waterloo, 200 University Avenue West,
5 Waterloo, Ontario, N2L 3G1, Canada

6 *Corresponding author: Email: janusz@uwaterloo.ca Phone: +1-519-888-4641. Fax: +1-
7 519-746-0435

8

ACCEPTED MANUSCRIPT

9

10 Abstracts

11 In LC-MS based untargeted analysis, data is collected at the peak or ion level, although
12 the investigated biochemistry processes occur at the compound or reaction level. To this
13 end, the presence of redundancy peaks such as co-eluted peaks, multi-chargers, adducts,
14 neutral loss, isotopologues, and fragments ions often muddle subsequent statistical data
15 analysis. In order to fill this gap, between peaks and compounds/reactions, independent
16 components must first be found at the peak level, then evaluated at the compound or
17 reaction levels. Based on paired mass distances (PMD), the algorithm GlobalStd, based
18 on retention time hierarchical cluster analysis and global analysis of PMDs within
19 clusters, is here proposed to extract independent peaks from raw LC-MS data. Following
20 its application, a structure/reaction directed analysis can then be used to evaluate
21 compounds at the structure or biochemistry reaction level, based on similar PMDs among
22 different retention times clusters. As a proof-of-concept, the developed statistical method
23 was applied to data obtained for *in vivo* SPME sampling on fish. In total, 277
24 independent peaks were demonstrated to stand for most of the variances found for the
25 total 1459 ions detected via LC-MS. Following, both known homologous series or
26 biological reactions along with unknown bio-processes, which may involve
27 oxidation/reduction reactions or homologous series, were analyzed via a
28 structure/reaction directed analysis. The findings of this analysis yielded interesting
29 information regarding the data, for instance denoting the possible occurrence of a
30 biosynthesis process involving L-Carnitine and its precursor 4-
31 Trimethylammoniobutanoic acid. Such PMD relationships could also aid in the screening
32 of annotation results. To this end, semi-quantitative analysis based on structure/reaction
33 directed analysis is also here proposed for further investigation of unknown patterns or
34 for removal of contaminants in metabolomics studies. The developed data-driven
35 algorithm has been included in a PMD package with a GUI interactive document, and is
36 freely available online (<https://github.com/yufree/pmd>).

37

38 **Keywords:** metabolomics; LC-MS; *in vivo*; SPME; algorithm

39

40

41

42

43 **1 Introduction**

44 Untargeted analysis based on Liquid chromatography–mass spectrometry (LC-MS)
45 has been applied in metabolomics,[\[1,2\]](#) petroleomics,[\[3\]](#) and environmental
46 analysis[\[4\]](#) for discoveries of unknown compounds associated with certain biotic
47 or abiotic chemical reactions. In such applications, compounds are usually first
48 analyzed in mass spectrometry as charged ions, while most of the downstream
49 analyses, such as group-wise differences, pathway analysis, and annotations, are
50 performed at the compound level, or between compounds.[\[5\]](#) In such cases,
51 charged ions' profiles can be further converted into peaks' profiles via peak
52 detection, using for instance feature detection algorithms such as centWave.[\[6\]](#)
53 However, peaks found from untargeted LC-MS analysis at a given retention time
54 are always comprised of a mixture of known and unknown co-eluted peaks, doubly
55 charged ions, adducts, neutral loss, isotopologues, fragments, or molecular ions[\[7\]](#).

56
57 The resulting ambiguity between the found peaks and their corresponding
58 compounds would thus affect any subsequent statistical analysis, as peak
59 intensities from any given compound would be in proportion to each other and
60 show a strong linear relation.[\[8,9\]](#) For example, in an analysis that yields 1000
61 peaks with statistical significant differences, such peaks may only correspond to
62 200 compounds. In such a case, then only twenty percent of the found peaks,
63 containing all pertinent information regarding the compounds of interest, would be
64 necessary for subsequent analysis. Also, peaks stemming from the same compound
65 would yield different sensitivities on LC-MS due to different intensities and signal-
66 to-noise ratios and introduce more uncertainties. Such peaks would also necessitate
67 validation with respect to potential false positive annotations of adducts or

68 fragmental ions.^[9] Besides, owing to different ionization processes or their
69 elemental composition, some compounds might only yield one peak, while others
70 may present multiple peaks with different kinds of adducts or isotopologues. This
71 redundancy in peaks would result in statistical bias once analysis such as multiple
72 comparisons^[10] with false discovery rate (FDR) control are carried out. Similar
73 issues have been discussed with respect to genomics research regarding genes with
74 dependence, as such occurrences have been shown to increase the variance of FDR
75 estimators.^[11–13]

76
77 Targeted qualitative and quantitative analyses usually involve the use of
78 standards,^[14] where for any given compound, qualitative and quantitative ions are
79 selected based on its standard's mass spectrum and/or retention time, with the
80 provision that selected ions for different compounds are distinctive from each other.
81 Following this rationale, in untargeted analysis, if only a few independent peaks
82 were to be selected among all of the identified peaks to stand for compounds prior
83 to further statistical analysis, then any subsequent statistical analyses could then be
84 performed at the compound level, as is performed in targeted analysis, thus
85 significantly simplifying data analysis and reducing uncertainty. Aiming to find
86 such peaks or remove redundancy peaks, previous studies have attempted to screen
87 mass-to-charge ratios (m/z) with the use of predefined adducts, neutral loss, or
88 chemical contaminants.^[9] However, unknown adducts, neutral loss, or chemical
89 contaminants may also play an important role in the profiles of untargeted peaks.^[7]
90 As the importance of such peaks cannot be reasonably predicted and must be
91 evaluated on a case-by-case basis, predefined rules based methods may thus miss
92 peaks from unknown background ions or adducts. Peak shapes^[15] or peak
93 intensities^[16–18] have been employed to identify pseudospectra of independent
94 compounds as a means to make further annotations or identifications. However, if

95 the purpose of research is to elucidate overall changes at the structural or reaction
96 level, then it can be reasonably assumed that one peak from each independent
97 compound would provide enough information.

98
99 Supposing that an independent peak for each compound is detected and selected,
100 compound identification would nonetheless still require the execution of additional
101 steps, such as tandem mass spectrometry analysis or database queries.**[16,18]** To
102 this end, specific tools have been developed to allow for predictions based on
103 MS/MS database data, such as the Global Natural Products Social Molecular
104 Networking (GNPS)**[19]** and Metlin.**[20]** Mass defect analysis, on the other hand,
105 can be employed to reveal unknown compounds**[14,21]** and compounds with same
106 sub-structures would show similar mass defect values. In petroleomics or
107 environmental analysis, extensions of the concepts of Kendrick mass defect
108 analysis have been employed to find homologous series**[22]** compounds with
109 different base units, such as $-\text{CH}_2-$, $-\text{O}-$, $-\text{CH}_2\text{O}-$, from high resolution mass
110 spectrometry data.**[21–23]** Besides, mass defect values could also be used to filter
111 drug metabolites from high resolution mass spectrometry data, since the
112 metabolites of certain parent compounds would show mass defect values within
113 50mDa of their parents compounds.**[24]**

114
115 We could extend the concept of mass defect to paired mass distance (PMD; the
116 distance between two mass-to-charge ratios), since chemical reactions would also
117 involve unique defect values between reactants. For instance, in environmental
118 analysis, a PMD of 33.96102Da is often used to screen for halogenated
119 contaminants, as this PMD corresponds to a dechlorination reaction that involves
120 an exchange between a hydrogen atom and a chlorine atom ($\text{H} \leftrightarrow \text{Cl}$)**[25]**. Thus,
121 PMD-based identification could be used to identify certain homologous series or

122 substitution reactions in biotic or abiotic processes. As such, further qualitative or
123 quantitative analysis could forego identifications of all detected compounds, only
124 focusing on compounds that present the same PMDs as a group, since such
125 compounds would have similar structures, or participate in the same chemical
126 reactions. Besides, it can be assumed that if a certain compound is involved in
127 multiple common PMDs, such compound would thus play an important role in the
128 untargeted profile of the found peaks.

129
130 However, employment of either methods, namely the identification of independent
131 peaks or PMD-based structure/reaction directed analysis, would necessitate a pre-
132 defined PMD for either adducts, neutral loss, isotopologues, or sub-structures as
133 shown in CSPP algorithm[18]. To this end, if a heuristic method could be
134 employed to find unknown PMDs based on the statistical properties of the LC-MS
135 peak profile, both known and unknown compounds belonging to adducts, neutral
136 loss, the same homologous series, or biochemistry reactions could thus be
137 identified. Once the LC-MS data is thus 'filtered', then subsequent semi-qualitative
138 or quantitative statistical analyses could be performed for those compounds as a
139 group, thus bypassing the need for identification of each peak found in the raw LC-
140 MS data. Further, such an approach to analysis might reveal unknown novel
141 mechanisms in untargeted studies, such as oxidation processes or substitution
142 reactions.

143
144 Selection of PMDs could include PMDs that correspond to certain structures or
145 reactions by element analysis. For instance, reactions involving Oxygen
146 (15.994915 amu), Phosphorus (30.973763 amu) and Sulfur (31.972072 amu)
147 would yield a PMD value corresponding to less than an integer number, such as
148 PMD 13.98 Da and 15.99 Da. On the other hand, reactions that involve Hydrogen

149 (1.007825 amu) and Nitrogen (14.003074 amu) atoms would always yield
150 reaction-related PMDs that are larger than an integer. For instance, for analysis of
151 data from samples collected from biological sources, PMD-based analysis could be
152 employed to infer which elements are involved in the specific metabolic reaction
153 under study. Such information would help point at the biochemical processes
154 associated with the studied phenomena without necessitating identification of each
155 detected feature.

156 For short-lived compounds, identification of reactions or dynamic changes
157 occurring among such compounds' structures would reveal important information
158 regarding their biological or environmental profile. While *in vivo* untargeted
159 studies can aid in the identification of previously unreported compounds,[26,27]
160 qualitative analysis of such "unknown" compounds captured by new analytical
161 methods remains a challenge. To this end, solid phase microextraction (SPME) has
162 been successfully applied towards analysis of *in vivo* biological processes to reveal
163 the presence of previously unreported short-lived compounds, which may have
164 gone undetected in analysis employing traditional sampling methods[26,28]. Thus,
165 *in vivo* SPME is presented as a suitable analytical platform to set up and validate a
166 statistical method for identification of unknown compounds.

167
168 In the current study, an algorithm, namely GlobalStd, is proposed to remove
169 redundancy peaks in LC-MS based non-targeted analysis, based on peaks' exact
170 mass and retention times. Following application of the algorithm, the resulting
171 independent peaks can then be submitted to a structure/reaction directed analysis at
172 the compound or reaction level. Such a method is designed to detect both known
173 and unknown compounds, as well as reaction relationship among compounds. As a
174 proof-of-concept, the developed method was employed towards untargeted
175 analysis of fish tissue via *in vivo* SPME.

176 **2 Materials and methods**

177

178 **2.1 Chemicals**

179

180 LC-MS grade acetonitrile, methanol, and water were purchased from Fisher
181 Scientific (Ottawa, ON, Canada). Hexane and acetone were purchased from
182 Sigma-Aldrich (Oakville, ON, Canada). Biocompatible SPME mixed mode probes
183 (45 μm thickness, 15mm length of coating) were provided by Supelco (Bellefonte,
184 PA, USA). Standards, including diazepam, nordiazepam, oxazepam, flunitrazepam,
185 lorazepam, testosterone, and progesterone were used as instrumental QC samples,
186 and purchased from Sigma-Aldrich (Oakville, ON, Canada).

187

188 **2.2 *In vivo* SPME sampling**

189

190 All experimental protocols were approved by and carried out in accordance with
191 guidelines established by the University of Waterloo Animal Care Committee
192 (AUPP #10–17). Rainbow trout (*Oncorhynchus mykiss*) (n=3) were purchased
193 from Silver Creek Aquaculture (Erin, ON Canada). Fish were acclimatized to
194 laboratory conditions for two weeks in non-chlorinated water. All fibers used in *in*
195 *vivo* sampling were preconditioned in methanol/water (50/50, v/v) prior to use.
196 Three fibers were used to sample each fish, and a total of three fish were sampled.

197

198 *In vivo* sampling of fish muscle tissue was carried out by inserting mixed mode
199 SPME fibers into the dorsal-epaxial muscle (near the dorsal fin) of fish after they
200 were anaesthetized with Tricaine mesylate and affixed to a foam bed. After
201 insertion of fibers, fish were allowed to recover in a bucket for a 20 minute period
202 while *in vivo* extractions were carried out. Once the extraction period was
203 concluded, fibers were pulled out, wiped with Kimwipes, and vortexed at 1500
204 rpm for 5s with ultrapure water to remove any matrix components on the fiber.

205 Desorption of fibers was performed with 300 μ L of acetonitrile/water (80/20, v/v)
206 as solvent for 90 min at 1,000 rpm vortex agitation. Extract solutions were
207 collected for instrumental analysis.

208

209 **2.3 Instrumental analysis**

210

211 An ACQUITY ultra performance liquid chromatography (UPLC) M-Class (UPLC)
212 instrument coupled with a mass spectrometer (Xevo G2-S QT of mass
213 spectrometer equipped with ZSpray TM ESI source) was used for instrumental
214 analysis of samples. Chromatography columns (Kinetex 1.7 μ m PFP, 100A, 100 x
215 2.1 mm) were eluted with mobile phase A (water with 0.1% Formic acid) and
216 mobile phase B (Acetonitrile with 0.1% Formic acid) at 80 μ L/min. The column
217 temperature was set at 30°C, and the samples were kept at 5°C. The injection
218 volume was 10 μ L. Gradient elution was as follows: 90% A was run for 1 min,
219 reduced to 10% in the following 7 min, then kept for 4 min. Following, mobile
220 phase A was increased back to 90% within 2 min, then maintained for 4 min until
221 the next injection.

222

223 The mass spectrometer was run in positive mode with spray voltage 3000V, cone
224 voltage 40V, and source offset 80V. The source temperature was 120°C and the
225 desolvation temperature was 350°C, with desolvation gas flow at 800L/h.
226 Acquisition mode was set as full scan mode, with a mass range of m/z 100-1000.
227 LockMass acquisition was employed to calibrate the mass spectrum, with a scan
228 time of 0.1s, an interval of 120s, and with Leucine enkephalin used as reference
229 material. Pool QC sample, and instrument QC samples were injected before and
230 after nine samples, blank solvent and blank fiber to assess the stability of the mass
231 spectrum throughout analysis. Such quality control showed a stable performance
232 during the analysis.

233

234 **2.4 Data processing**

235

236 Following instrumental analysis, raw data was exported from the instrument and
237 converted into mzxml format for further data analysis. Once optimized parameters
238 were attained via employment of the IPO package[29] on pool QC samples,
239 XCMS[30] was used to extract peaks. The GlobalStd algorithm was then
240 employed in a structure/reaction directed analysis to evaluate the obtained profiles
241 from fish *in vivo* sampling. Metlin was used to tentatively annotate the peaks and
242 obtain chemical names via comparisons of chemical formulas, with an accuracy of
243 less than 5 ppm. Annotation was employed to validate the results of the
244 structure/reaction directed analysis.

245

246 **2.5 GlobalStd algorithm**

247

248 **STEP 1: Retention time cluster analysis**

249 The algorithm GlobalStd was developed to find independent peaks from peak
250 retention time and mass-to-charge ratio profiles. As shown in scheme 1, the first
251 step of GlobalStd encompasses the aggregation of peak groups based on a retention
252 time hierarchical cluster analysis.[31] Such groups include components separated
253 by chromatography that are relatively independent of each other. Once this
254 analysis is concluded, then PMD analysis can be used to screen potential
255 redundancy peaks.

256

257 **STEP 2: Paired mass distance (PMD) analysis**

258 Redundancy peaks from same compounds should be discarded for further
259 structure/reaction directed analysis. As shown in work by Mahieu et al.,⁷ unknown
260 adducts or background ions could be revealed by frequent intrinsic relationship
261 analysis. However, the presence of co-eluted compounds can make such

262 frequency-based methods exclude unknown redundancy peaks from unrelated
263 compounds. On the other hand, as doubly charged ions would show a PMD around
264 0.5, PMD analysis can enable the exclusion of these mass pairs from further
265 discussion. To avoid the inclusion of common isotopologues, e.g., peaks with ^{12}C
266 and ^{13}C , mass distance pairs around 1 and 2 would be treated as isotopologues
267 groups,[23] and any additional PMD analysis would only include isotopologues
268 with lower mass-to-charge ratios. As such, ions identified via PMD analysis will
269 not have isotopologues or doubly charged ions among the data carried forward for
270 further analysis.

271
272 Following the above discussed steps, further PMD analysis can then performed
273 based on the 'global' properties of the PMDs found in each retention time group. If
274 a specific PMD were to appear multiple times in different retention time groups,
275 then such PMD would be assumed to reflect universal paired relationships, such as
276 adducts or neutral loss. At the same time, most of the fragmented ions, co-eluted
277 compounds, or contaminated ions would be removed for further analysis, as their
278 PMDs are unlikely to appear in multiple retention time groups as compared with
279 adducts or neutral loss. Since only PMDs within the same retention time group are
280 addressed in this step, PMDs between independent compounds would thus not be
281 captured.

282

283 STEP 3: Selection of independent peaks

284 The workflow of Step 3 is illustrated in the right part of Scheme 1. Here, within
285 each retention time group, the remaining peaks are grouped into one of two
286 groupings: one that contains singles peak in the retention time group, and another
287 to encompass multiple peaks in the retention time group. Here, single peaks are
288 kept as independent components. The grouping with multiple peaks, on the other

289 hand, is further grouped into another two categories. One category is comprised of
290 peaks with isotopologue peaks, while the other is comprised of peaks with ‘global’
291 PMDs. For retention time groups that contain multiple groups, only the largest
292 mass-to-charge ratios are selected as potential molecular ions or base peaks.

293
294 For retention time groups with isotopologues or ‘global’ PMDs, we could further
295 divide their peaks into three parts: one with isotopologues peaks and no ‘global’
296 paired masses, one with ‘global’ paired masses and no isotopologues, and one last
297 group, containing both isotopologues and ‘global’ paired masses. For the first
298 group, since ^{12}C containing isotopologues often show higher intensities than
299 isotopologues with ^{13}C , smaller ions are then kept as independent ions. For the
300 second group, all ions with smaller mass-to-charge ratios in the ‘global’ paired
301 masses are treated as independent peaks. For the third group, all isotopologues
302 with lower mass-to-charge ratios are first extracted. Then, aiming to remove all
303 isotopologues adducts, the mass distances among the remaining isotopologues are
304 calculated, and only the lower mass isotopologues that appear in the ‘global’
305 PMDs are kept. Other isotopologue ions can also be kept as potential independent
306 peaks, even if they are not in the ‘Global’ PMDs.

307
308 Once all these steps are concluded, and most if not all repeated peaks,
309 isotopologues, and adduct-related peaks are removed, the peaks from all of the
310 above groups can be combined together as independent peaks, and carried forward
311 for further analytical analysis. While this step is aimed at removing isotopologues
312 and adduct related peaks, the remaining peaks could still contain some adducts ions
313 if these compounds are only shown as adduct ions.

314

315 In summary, the goal of the GlobalStd algorithm is to use a minimum amount of
316 peaks to stand for the significant amount of peaks generally found in untargeted
317 analyses by removing redundancy peaks from the same compounds prior to further
318 analytical analysis. To this end, the presented algorithm requires at least two
319 parameters: the cutoff of the retention time hierarchical cluster analysis, and the
320 bottom threshold number of retention time groups for global PMD searches. For
321 example, a threshold of 10 would mean that the selected PMD should appear in at
322 least 10 different retention time groups. Since we employed cluster analysis, the
323 resolution of the chromatography separation could be controlled by the cutoffs of
324 distances between retention time groups. Such a cutoff should reflect the
325 separation capacity of the employed chromatography columns. Selection of an
326 appropriate bottom threshold number for PMD searches, on the other hand, would
327 ensure that retention time groups for PMD analysis can be determined by explorer
328 analysis of the PMDs profiles so as to include all m/z with known PMDs.

329 **2.6 Structure/reaction directed analysis**

330
331 PMDs can also be used to group compounds in structure/reaction directed analyses.
332 Here PMDs for peaks in different retention time groups are used instead of PMDs
333 of the same retention time groups, as is the case for the GlobalStd algorithm
334 application. To this end, such PMDs would not indicate adducts or background
335 ions, since those peaks are supposedly coming from different compounds. These
336 PMDs may nonetheless be related to certain homologous series or chemistry
337 reactions. To this end, a frequency cutoff could be set to investigate universal
338 homologous series or chemistry reaction related compounds. The presence of
339 isomers would increase the frequencies of certain mass-to-charge ratios, thus
340 ensuring that only one of each isomer remains in the data carried forward for
341 frequency analysis. Such structure/reaction directed analysis could be performed at

342 the peak level without employing the GlobalStd algorithm. However, peaks
343 stemming from the same compounds would be cumbered by additional noise in the
344 frequency, as shown in the following section.

345
346 To make it clear, GlobalStd algorithm is different from published methods like
347 DeltaMS[32] or MSClust[33]. For DeltaMS, mass distances are used to find
348 isotopologues relationship[32]. However, our methods also used such relationship
349 to find adducts, neutral losses, homologous series or chemistry reactions.
350 MSClust[33] use intensity-based cluster analysis to reduce the peaks into
351 compounds while PMD method only use paired mass distances. As we will show in
352 the demonstrated data, our method could show a similar result compared with
353 intensity-based methods. However, since intensity was not used to find
354 independent peaks, our method is robust for the uncertainty in intensity
355 measurement. Another important difference is that our method doesn't use pre-
356 defined neutral losses, adducts lists, homologous series or reaction. All the findings
357 are based on relationship frequency in the data and only the high frequency paired
358 mass distance relationships are kept for further investigation. Current methods such
359 as mass defect, or could not find unknown reactions or adducts while our methods
360 could reveal them if they show a highly frequency in the peaks profile. As for
361 structure/reaction directed analysis, similar but totally different way has been used
362 to find metabolites for known compounds[34]. However, our method directly uses
363 the frequency of paired mass relationships to screen and reveal both known and
364 unknown structures or reactions.

365
366 Both the GlobalStd algorithm and the structure/reaction directed analysis workflow
367 have been included in the PMD package, which is freely available online
368 (<https://github.com/yufree/pmd>). All the documents for this package could be

369 found online (<https://yufree.github.io/pmd/>). A graphical user interface (GUI) to
370 perform the presented PMD-based methods was also included in this software
371 package as interactive documents. Experimental data from *in vivo* SPME sampling
372 are also attached in this package for reproducible research purposes.

373

ACCEPTED MANUSCRIPT

374
375

376 **3 Results and Discussion**

377
378
379

3.1 Retention time groups

380 A total of 1459 peaks were extracted from *in vivo* SPME samples across 9 samples.
381 As shown in Figure S1, 75 retention time groups were found in *in vivo* SPME
382 datasets. Under the employed chromatography conditions, hydrophilic compounds
383 eluted first, followed by lipophilic compounds. Indeed, some hydrophilic
384 compounds were observed to not retain on the employed column, and to co-elute at
385 the very beginning (see retention time group 6). Some patterns, such as
386 homologous series, could also be observed in the raw data as such compounds
387 eluted sequentially, with an increasing mass-to-charge ratio (m/z). However, the
388 majority of peaks formed what appeared to be a random pattern on the retention
389 time- m/z profile. Peaks within certain retention time groups could be either co-
390 eluted compounds or peaks from same compounds. Hierarchical cluster analysis
391 separates those peaks with a cutoff of 10, which means the complete linkage
392 distances between each retention time group is larger than 10s.

393

394 In summary, retention time hierarchical cluster analysis could aid in the search for
395 relatively independent fractions. Following, PMD-based filtering could be applied
396 within each retention time group to further reduce peaks into potential independent
397 peaks.

398

3.2 PMD analysis

399

400 A PMD analysis with cutoff of 10 for the frequency of PMDs between RT cluster
401 for independent peaks of the *in vivo* data indicated 8 retention groups with single
402 peaks. Additionally, 631 isotopologue-related paired mass and 685 multi-charger
403

404 related paired mass peaks were found. As shown in figure 1, among the PMDs to
405 appear in more than 10 retention time groups, 10 unique PMDs (which retain 2
406 digits after the decimal point), involving 431 peaks, were kept out of 443 paired
407 mass peaks. Some PMDs were treated as adducts (such as 21.98Da for adducts
408 between H^+ and Na^+ , and 17.03Da for adducts between H^+ and NH_4^+)[35] while
409 some were treated as neutral loss (such as 18.01Da for H_2O).[36] Some polymer-
410 related PMDs, such as PMD 28.03Da ($-C_2H_4-$) and 44.03Da ($-C_2H_4O-$) were also
411 found. Unknown adducts/neutral loss such as PMD 45.06Da ($-C_2H_7N-$) and
412 66.01Da ($-C_4H_2O-$ or $-C_2N_3-$) were also identified in the presently discussed data.
413 Conversely, PMD 23.0760Da, shown in Mahieu et. al's work,[7] was not found in
414 this dataset, which means such a PMD may be related to case-by-case unknown
415 background ions. It should be also noted that Mahieu et. al's analysis directly use
416 global paired mass distances for all mass pairs, while our algorithm only employs
417 the mass distances within each retention time group. Besides, the median Pearson
418 correlation coefficients of the PMD's intensity is 0.88, which implies these paired
419 peaks stem from the same compounds. In summary, PMD analysis within retention
420 time groups could show both known and unknown adducts or background ions
421 from the m/z - retention time profile.

422 423 3.3 Independent peaks selection

424 Application of the GlobalStd algorithm on the data yielded 277 independent ions.
425 As shown in figure S2, ions found by GlobalStd could fit into different scenarios,
426 including groups with lots of co-eluted peaks. Since the developed algorithm only
427 uses m/z and retention times, intensities can thus be further applied to validate the
428 selections. Likewise, Principal component analysis (PCA)[37] can be used to
429 assess changes between score plots of the raw data (containing all peaks) versus
430 that of the selected peaks. As shown in figure 2, PCA similarities would indicate
431

432 that the selected peaks, representing around 20% of the original peaks, sufficiently
433 capture variances from all peaks. Considering that this analysis only employs peak
434 mass and retention times, a correlation analysis based on peak intensities could be
435 used as an independent test to further screen peaks.

436
437 Further validation can be carried out based on a statistical analysis of peak
438 intensities. As three fish were sampled with three SPME fibers each, no statistical
439 differences should be found among biological replicates. From the raw peaks, 86
440 peaks out of 1459 peaks showed statistical differences among three fish (F test, p-
441 value cutoff 0.05). With a p-value cutoff of 0.05 for multiple comparisons, 73
442 (1459×0.05) peaks with significant statistical differences were identified as false
443 positives. From the independent peaks, 17 peaks out of 277 peaks were found to
444 yield statistical differences (F test, p-value cutoff 0.05). Of these, 14 (277×0.05)
445 peaks should be identified as false positives, with p-value cutoff 0.05 for multiple
446 comparisons in independent peaks. After applying a false discovery rate control (q-
447 value cutoff 0.05), no peaks could be identified as true discoveries, in either the
448 raw or selected data. Thus, the statistical analysis would indicate that the
449 algorithm-selected peaks retain information similar to that in the raw data.
450 However, the validation of those peaks are reduced a lot by focused on
451 independent peaks.

452
453 In summary, the GlobalStd algorithm can be used to reduce peak numbers with
454 minimal loss of information. As a next step, the peaks selected by the algorithm
455 can then be submitted to structure/reaction directed analysis.

456 **3.4 Structure/reaction directed analysis**

457

458 277 selected independent peaks were imported for structure directed analysis. Here,
459 only PMDs among different retention time groups were considered for the selected
460 peaks. This setting forced the structure/reaction directed analysis to use peaks
461 which could be separated by chromatography. In total, 19 PMD groups with a
462 frequency larger than 10, as shown in figure 3, were found. All 277 peaks were
463 submitted to Metlin for their chemical formula, with accuracy setting of less than 5
464 ppm. The settings $[M+H]^+$, $[M+Na]^+$, $[M+NH_4]^+$ and $[M-H_2O+H]^+$ were selected
465 for the database search, according to the PMD analysis above. Potential
466 structure/reactions were then directly investigated via comparisons of chemical
467 formula. As shown in table S1, 119 peaks were involved in those 19 PMD groups.
468 This would indicate tentative identification of a variety of compounds involved in
469 networks of multiple chemical reactions in certain biotic or abiotic systems. For
470 example, $C_{22}H_{41}NO_2$ (m/z 352.3214), which appeared in 7 different PMD groups,
471 was tentatively identified as anandamide, a reported active compound in living
472 systems.[38]

473
474 Some of the identified PMD groups highlighted in this analysis have been already
475 associated with known structures or certain bio-processes. For instance, a PMD of
476 0 indicates isomers, while a PMD of 13.98Da could indicate the exchange of an
477 oxygen atom for two hydrogen atoms, which is associated to an oxidation process
478 followed by H_2O elimination.[35,39] For example, $C_{24}H_{36}O_5$ (m/z 405.2616) and
479 $C_{24}H_{38}O_4$ (m/z 391.2835)'s ions were tentatively identified with accuracy less than
480 5 ppm, and a statistically significant intensity correlation (pearson correlation
481 coefficient, 0.8427). Such a relationship might denote the presence of an oxidation
482 process of 3β -Hydroxy-6-oxo-5 α -cholan-24-oic acid, according to tentative
483 annotation from Metlin.[20] Likewise, a PMD of 15.99Da might indicate the
484 addition or removal of oxygen atoms, or an oxidation process.^{31,35} For example,

485 $C_7H_{15}NO_2$ (m/z 162.1128) and $C_7H_{15}NO_3$ (m/z 146.1183) could be L-Carnitine and
486 its precursor 4-Trimethylammoniobutanoic acid, as supported by their intensity
487 correlation (pearson correlation coefficient, 0.9519). This biosynthesis process has
488 been reported to occur in humans,[40] and may also similarly occur in fish, since
489 L-Carnitine is also found in fish.[41]

490
491 Some of the acquired PMD values could be related to homologous series such as
492 PMD 14.02Da, 28.03Da, and 58.04Da. These in turn could be related to
493 substructures of $-CH_2-$, $-C_2H_4-$ and $-C_3H_6O-$, respectively. Such substructures
494 could be found in fatty acids^{31,35} or surfactant[21]. As shown in figure 4, a series of
495 seven compounds, from m/z 425.3120 to m/z 773.5662, and with a PMD of
496 58.04Da, were identified in the data. The chromatograph also showed a linear
497 elution process, with regular increasing distances. However, a Metlin[20]search
498 failed to yield corresponding compounds. Previous works^{31,35} have treated PMD
499 58.04Da as acetone condensation, although such a process might not occur
500 between compounds. The identified PMD might also be related to polymers such
501 as Polypropylene Glycol, since their mass spectrum covers the peaks found in our
502 research.[42] While these peaks were not found in fiber control, they may very
503 well be contaminants or unknown compounds.

504
505 Some unknown PMDs might need further validation analysis. For example, a PMD
506 16.03 could be related to a mass difference of one carbon atom and four hydrogen
507 atoms. On the other hand, this PMD may correspond to a combination of removing
508 the substructure $-CH_2-$ and a dehydrogenation reaction. In *in vivo* SPME sampling,
509 such a PMD was found between $C_{24}H_{50}NO_7P$ (m/z 496.3410) and $C_{23}H_{46}NO_7P$
510 (m/z 480.3100) (pearson intensity correlation coefficient of 0.9456), which might
511 indicate a conversion between two kinds of phosphatidylcholine(PC) ($C_{24}H_{50}NO_7P$

512 \leftrightarrow $C_{23}H_{46}NO_7P$). Another possibility however, is that both of these two
513 compounds stem from the same parent compounds (unknown parent compounds \rightarrow
514 $C_{24}H_{50}NO_7P + C_{23}H_{46}NO_7P$). Both types of reactions would show a high
515 correlation coefficient between the two compounds. If this PMD were to appear
516 with high frequency for certain metabolites across the independent peaks, then it
517 could be reasonably concluded that such reactions are not the result of a random
518 combination of two compounds. While all the annotations made in this work need
519 further validation, such as MS/MS analysis or data-based predictions, some
520 preliminary conclusions can be nonetheless drawn at the chemical formula level
521 for the unknown parts of this non-targeted analysis.

522
523 Compounds from homologous series or similar biochemistry reactions, such as
524 lipids, might show response factors with regularity on mass spectrometry
525 analysis.^[43] The average responses from certain mass defect groups could be used
526 for a semi-quantitative evaluation of those unknown homologous or reactions in
527 samples. Figure S3 shows the relative standard deviations (RSD%) of compounds
528 in each group among the three fish. The peaks can be further filtered for certain
529 homologous series or similar biochemistry reactions by assigning a threshold based
530 on the attained RSD%. If a given PMD group shows significant average intensity
531 changes among the two conditions, then this change can be directly used to
532 quantitate certain homologous series or biochemistry reactions, which would allow
533 for a circumvention of the use of standards to validate these compounds. Further,
534 an established linear relationship between paired masses could be also used to filter
535 reasonable peaks for subsequent semi-quantitative structure or reaction analysis.

536
537 To this end, such an analysis was also performed on raw peak data, without prior
538 application of the GlobalStd algorithm. As shown in Figure S4, PMD values

539 yielded a chaotic distribution, with peaks from same compounds, and much noise
540 in the frequency. While known adducts could be used to filter data in each PMD
541 group, unknown adduct ions, such as PMD 66.01, would still be present in the data.
542 Thus, application of the GlobalStd algorithm would be necessary to remove both
543 known and unknown peaks from the same compounds.

544

545 **4 Conclusions**

546

547 The current work proposes a data-driven method to evaluate untargeted data at the
548 compound, homologous series, or biochemistry reaction levels without the use of
549 standards or intensity data. The presented methods could be used to remove
550 redundancy peaks from data profile and to select independent peaks for further
551 structure/reaction directed analysis. As this process is software automated and
552 based on a heuristic search, it enables the unveiling of both known and unknown
553 relationships between the peaks. PMD values can be used to elucidate bio-
554 processes at the reaction level, as well as to aid in more accurate peak annotations
555 based on their oxidation-reduction properties. To this end, the establishment of a
556 database of PMDs and their corresponding homologous and reactions might aid in
557 much easier exploration of “unknown unknown” compounds.

558

559

560 **Acknowledgments**

561
562 We gratefully acknowledge Professor Mark R. Servos and Ms. Leslie Bragg for
563 kindly providing us with fish and the use of a facility for *in vivo* sampling as well
564 as for their assistance in the sampling. We are truly grateful for Ms. Sofia Lendor'
565 help on instrumental analysis of the samples.

566
567 **Funding**

568
569 This research was financially supported by the Industrial Research Chair of the
570 National Sciences and Engineering Research Council of Canada (NSERC-IRC).
571 We are also very grateful to Waters Corporation for kindly providing the seed
572 instrument Xevo G2-S QToF to our laboratory. The authors thankfully
573 acknowledge Supelco (Millipore Sigma) for providing the mixed-mode SPME
574 fibers used in this study.

575

576

577 **References**

- 578 [1] A.C. Schrimpe-Rutledge, S.G. Codreanu, S.D. Sherrod, J.A. McLean,
579 Untargeted Metabolomics Strategies-Challenges and Emerging Directions, *J. Am.*
580 *Soc. Mass Spectrom.* 27 (2016) 1897–1905.
- 581 [2] A. Scalbert, L. Brennan, O. Fiehn, T. Hankemeier, B.S. Kristal, B. van
582 Ommen, E. Pujos-Guillot, E. Verheij, D. Wishart, S. Wopereis, Mass-
583 spectrometry-based metabolomics: limitations and recommendations for future
584 progress with particular focus on nutrition research, *Metabolomics.* 5 (2009) 435–
585 458.
- 586 [3] N.E. Oro, C.A. Lucy, Analysis of the Nitrogen Content of Distillate Cut Gas
587 Oils and Treated Heavy Gas Oils Using Normal Phase HPLC, Fraction Collection
588 and Petroleomic FT-ICR MS Data, *Energy Fuels.* 27 (2013) 35–45.
- 589 [4] S. Samanipour, M. Reid, K. Bæk, K.V. Thomas, Combining a deconvolution
590 and a universal library search algorithm for the non-target analysis of data
591 independent LC-HRMS spectra, *Environ. Sci. Technol.* (2018).
592 doi:10.1021/acs.est.8b00259.
- 593 [5] A. Alonso, S. Marsal, A. Julià, Analytical methods in untargeted
594 metabolomics: state of the art in 2015, *Front Bioeng Biotechnol.* 3 (2015) 23.
- 595 [6] R. Tautenhahn, C. Böttcher, S. Neumann, Highly sensitive feature detection
596 for high resolution LC/MS, *BMC Bioinformatics.* 9 (2008) 504.
- 597 [7] N.G. Mahieu, G.J. Patti, Systems-Level Annotation of a Metabolomics Data
598 Set Reduces 25 000 Features to Fewer than 1000 Unique Metabolites, *Anal. Chem.*
599 89 (2017) 10397–10406.

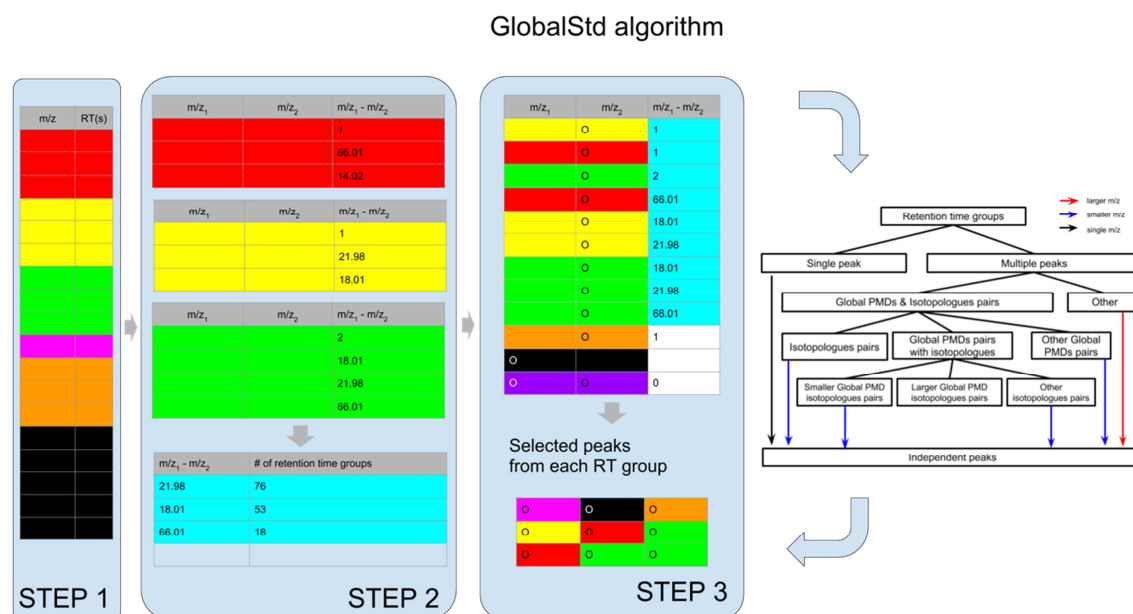
- 600 [8] John Wiley & Sons, Inc., ed., *Mass Spectrometry*, in: Kirk-Othmer
601 *Encyclopedia of Chemical Technology*, John Wiley & Sons, Inc., Hoboken, NJ,
602 USA, 2000: p. 542.
- 603 [9] X. Domingo-Almenara, J.R. Montenegro-Burke, H.P. Benton, G. Siuzdak,
604 *Annotation: A Computational Solution for Streamlining Metabolomics Analysis*,
605 *Anal. Chem.* 90 (2018) 480–489.
- 606 [10] J. Xia, D.I. Broadhurst, M. Wilson, D.S. Wishart, *Translational biomarker*
607 *discovery in clinical metabolomics: an introductory tutorial*, *Metabolomics*. 9
608 (2013) 280–299.
- 609 [11] A.B. Owen, *Variance of the Number of False Discoveries*, *J. R. Stat. Soc.*
610 *Series B Stat. Methodol.* 67 (2005) 411–426.
- 611 [12] J.T. Leek, J.D. Storey, *Capturing heterogeneity in gene expression studies*
612 *by surrogate variable analysis*, *PLoS Genet.* 3 (2007) 1724–1735.
- 613 [13] Y. Benjamini, D. Yekutieli, *The Control of the False Discovery Rate in*
614 *Multiple Testing under Dependency*, *Ann. Stat.* 29 (2001) 1165–1188.
- 615 [14] L.D. Roberts, A.L. Souza, R.E. Gerszten, C.B. Clish, *Targeted*
616 *metabolomics*, *Curr. Protoc. Mol. Biol.* Chapter 30 (2012) Unit 30.2.1–24.
- 617 [15] B. Zhou, J.F. Xiao, L. Tuli, H.W. Ransom, *LC-MS-based metabolomics*,
618 *Mol. Biosyst.* 8 (2012) 470–481.
- 619 [16] K. Uppal, D.I. Walker, D.P. Jones, *xMSannotator: An R Package for*
620 *Network-Based Annotation of High-Resolution Metabolomics Data*, *Anal. Chem.*
621 89 (2017) 1063–1067.
- 622 [17] C.D. Broeckling, F.A. Afsar, S. Neumann, A. Ben-Hur, J.E. Prenni,
623 *RAMClust: a novel feature clustering method enables spectral-matching-based*
624 *annotation for metabolomics data*, *Anal. Chem.* 86 (2014) 6812–6817.
- 625 [18] K. Morreel, Y. Saeys, O. Dima, F. Lu, Y. Van de Peer, R. Vanholme, J.
626 *Ralph, B. Vanholme, W. Boerjan, Systematic structural characterization of*

- 627 metabolites in *Arabidopsis* via candidate substrate-product pair networks, *Plant*
628 *Cell*. 26 (2014) 929–945.
- 629 [19] M. Wang, J.J. Carver, V.V. Phelan, L.M. Sanchez, N. Garg, Y. Peng, D.D.
630 Nguyen, J. Watrous, C.A. Kapon, T. Luzzatto-Knaan, C. Porto, A. Bouslimani,
631 A.V. Melnik, M.J. Meehan, W.-T. Liu, M. Crüsemann, P.D. Boudreau, E.
632 Esquenazi, M. Sandoval-Calderón, R.D. Kersten, L.A. Pace, R.A. Quinn, K.R.
633 Duncan, C.-C. Hsu, D.J. Floros, R.G. Gavilan, K. Kleigrew, T. Northen, R.J.
634 Dutton, D. Parrot, E.E. Carlson, B. Aigle, C.F. Michelsen, L. Jelsbak, C.
635 Sohlenkamp, P. Pevzner, A. Edlund, J. McLean, J. Piel, B.T. Murphy, L. Gerwick,
636 C.-C. Liaw, Y.-L. Yang, H.-U. Humpf, M. Maansson, R.A. Keyzers, A.C. Sims,
637 A.R. Johnson, A.M. Sidebottom, B.E. Sedio, A. Klitgaard, C.B. Larson, C.A.B. P,
638 D. Torres-Mendoza, D.J. Gonzalez, D.B. Silva, L.M. Marques, D.P. Demarque, E.
639 Pociute, E.C. O’Neill, E. Briand, E.J.N. Helfrich, E.A. Granatosky, E. Glukhov, F.
640 Ryffel, H. Houson, H. Mohimani, J.J. Kharbush, Y. Zeng, J.A. Vorholt, K.L.
641 Kurita, P. Charusanti, K.L. McPhail, K.F. Nielsen, L. Vuong, M. Elfeki, M.F.
642 Traxler, N. Engene, N. Koyama, O.B. Vining, R. Baric, R.R. Silva, S.J. Mascuch,
643 S. Tomasi, S. Jenkins, V. Macherla, T. Hoffman, V. Agarwal, P.G. Williams, J.
644 Dai, R. Neupane, J. Gurr, A.M.C. Rodríguez, A. Lamsa, C. Zhang, K. Dorrestein,
645 B.M. Duggan, J. Almaliti, P.-M. Allard, P. Phapale, L.-F. Nothias, T. Alexandrov,
646 M. Litaudon, J.-L. Wolfender, J.E. Kyle, T.O. Metz, T. Peryea, D.-T. Nguyen, D.
647 VanLeer, P. Shinn, A. Jadhav, R. Müller, K.M. Waters, W. Shi, X. Liu, L. Zhang,
648 R. Knight, P.R. Jensen, B.O. Palsson, K. Pogliano, R.G. Linington, M. Gutiérrez,
649 N.P. Lopes, W.H. Gerwick, B.S. Moore, P.C. Dorrestein, N. Bandeira, Sharing and
650 community curation of mass spectrometry data with Global Natural Products
651 Social Molecular Networking, *Nat. Biotechnol.* 34 (2016) 828–837.
- 652 [20] C. Guijas, J.R. Montenegro-Burke, X. Domingo-Almenara, A. Palermo, B.
653 Warth, G. Hermann, G. Koellensperger, T. Huan, W. Uritboonthai, A.E. Aisporna,

- 654 D.W. Wolan, M.E. Spilker, H.P. Benton, G. Siuzdak, METLIN: A Technology
655 Platform for Identifying Knowns and Unknowns, *Anal. Chem.* 90 (2018) 3156–
656 3164.
- 657 [21] M. Loos, H. Singer, Nontargeted homologue series extraction from
658 hyphenated high resolution mass spectrometry data, *J. Cheminform.* 9 (2017) 12.
- 659 [22] L. Sleno, The use of mass defect in modern mass spectrometry, *J. Mass*
660 *Spectrom.* 47 (2012) 226–236.
- 661 [23] E.M. Thurman, I. Ferrer, The isotopic mass defect: a tool for limiting
662 molecular formulas by accurate mass, *Anal. Bioanal. Chem.* 397 (2010) 2807–
663 2816.
- 664 [24] H. Zhang, D. Zhang, K. Ray, M. Zhu, Mass defect filter technique and its
665 applications to drug metabolite identification by high-resolution mass spectrometry,
666 *J. Mass Spectrom.* 44 (2009) 999–1016.
- 667 [25] K.J. Jobst, L. Shen, E.J. Reiner, V.Y. Taguchi, P.A. Helm, R. McCrindle, S.
668 Backus, The use of mass defect plots for the identification of (novel) halogenated
669 contaminants in the environment, *Anal. Bioanal. Chem.* 405 (2013) 3289–3297.
- 670 [26] D. Vuckovic, I. de Lannoy, B. Gien, R.E. Shirey, L.M. Sidisky, S. Dutta, J.
671 Pawliszyn, In vivo solid-phase microextraction: capturing the elusive portion of
672 metabolome, *Angew. Chem. Int. Ed Engl.* 50 (2011) 5344–5348.
- 673 [27] O. Teahan, S. Gamble, E. Holmes, J. Waxman, J.K. Nicholson, C. Bevan,
674 H.C. Keun, Impact of analytical bias in metabonomic studies of human blood
675 serum and plasma, *Anal. Chem.* 78 (2006) 4307–4318.
- 676 [28] V. Bessonneau, J. Ings, M. McMaster, R. Smith, L. Bragg, M. Servos, J.
677 Pawliszyn, In vivo microsampling to capture the elusive exposome, *Sci. Rep.* 7
678 (2017) 44038.

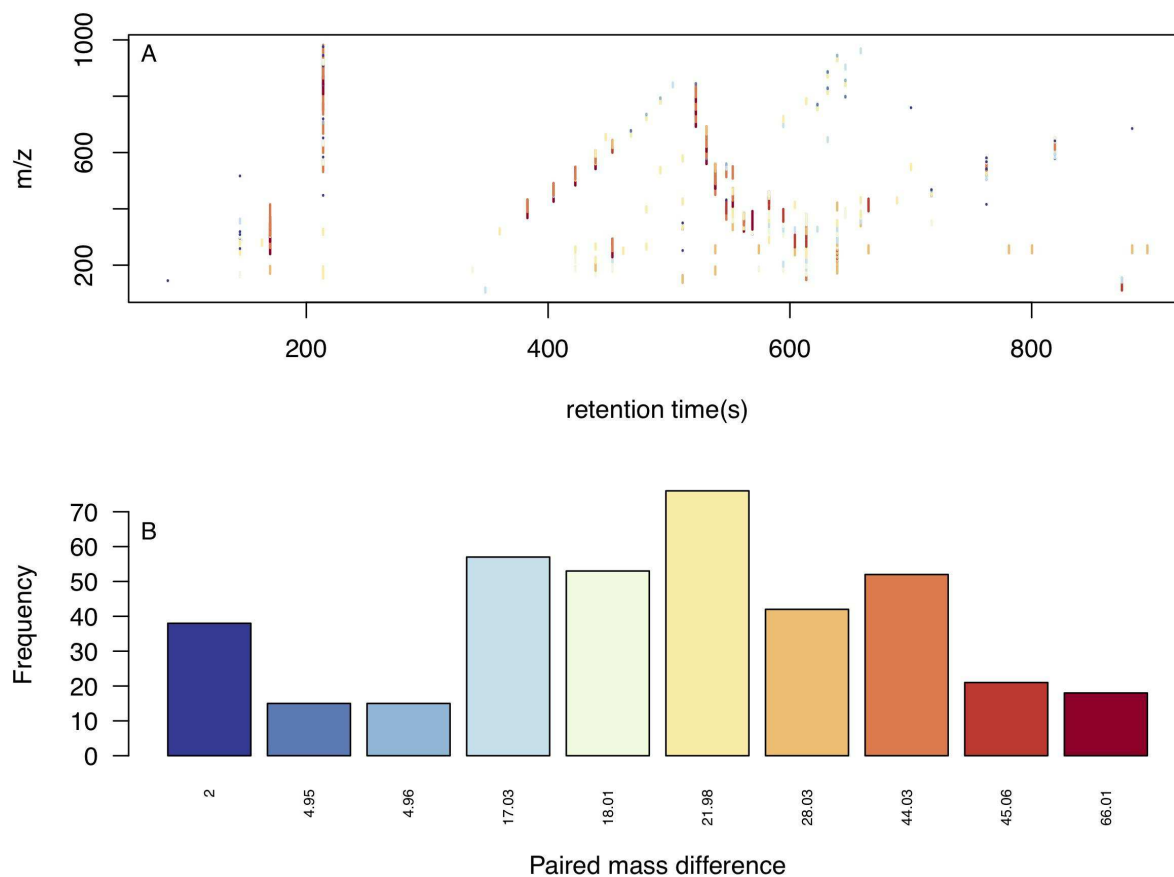
- 679 [29] G. Libiseller, M. Dvorzak, U. Kleb, E. Gander, T. Eisenberg, F. Madeo, S.
680 Neumann, G. Trausinger, F. Sinner, T. Pieber, C. Magnes, IPO: a tool for
681 automated optimization of XCMS parameters, *BMC Bioinformatics*. 16 (2015) 118.
- 682 [30] C.A. Smith, E.J. Want, G. O'Maille, R. Abagyan, G. Siuzdak, XCMS:
683 processing mass spectrometry data for metabolite profiling using nonlinear peak
684 alignment, matching, and identification, *Anal. Chem.* 78 (2006) 779–787.
- 685 [31] C.C. Bridges Jr., Hierarchical Cluster Analysis, *Psychol. Rep.* 18 (1966)
686 851–854.
- 687 [32] T.U.H. Baumeister, N. Ueberschaar, W. Schmidt-Heck, J.F. Mohr, M.
688 Deicke, T. Wichard, R. Guthke, G. Pohnert, DeltaMS: a tool to track isotopologues
689 in GC- and LC-MS data, *Metabolomics*. 14 (2018) 41.
- 690 [33] Y.M. Tikunov, S. Laptinok, R.D. Hall, A. Bovy, R.C.H. de Vos, MSClust: a
691 tool for unsupervised mass spectra extraction of chromatography-mass
692 spectrometry ion-wise aligned data, *Metabolomics*. 8 (2012) 714–718.
- 693 [34] E.A.P. Ekanayaka, M.D. Celiz, A.D. Jones, Relative mass defect filtering of
694 mass spectra: a path to discovery of plant specialized metabolites, *Plant Physiol.*
695 167 (2015) 1221–1232.
- 696 [35] B.O. Keller, J. Sui, A.B. Young, R.M. Whittal, Interferences and
697 contaminants encountered in modern mass spectrometry, *Anal. Chim. Acta.* 627
698 (2008) 71–81.
- 699 [36] R.M. Smith, ed., Neutral Losses and Ion Series, in: *Understanding Mass*
700 *Spectra: A Basic Approach*, John Wiley & Sons, Inc., Hoboken, NJ, USA, 2004:
701 pp. 121–149.
- 702 [37] G.H. Dunteman, *Principal Components Analysis*, SAGE, 1989.
- 703 [38] D.L. Lynch, P.H. Reggio, Cannabinoid CB1 receptor recognition of
704 endocannabinoids via the lipid bilayer: molecular dynamics simulations of CB1

- 705 transmembrane helix 6 and anandamide in a phospholipid bilayer, *J. Comput.*
706 *Aided Mol. Des.* 20 (2006) 495–509.
- 707 [39] H. Tong, D. Bell, K. Tabei, M.M. Siegel, Automated data massaging,
708 interpretation, and e-mailing modules for high throughput open access mass
709 spectrometry, *J. Am. Soc. Mass Spectrom.* 10 (1999) 1174–1187.
- 710 [40] C.J. Rebouche, A.G. Engel, Tissue distribution of carnitine biosynthetic
711 enzymes in man, *Biochimica et Biophysica Acta (BBA) - General Subjects.* 630
712 (1980) 22–29.
- 713 [41] C. Rigault, F. Mazué, A. Bernard, J. Demarquoy, F. Le Borgne, Changes in
714 l-carnitine content of fish and meat during domestic cooking, *Meat Sci.* 78 (2008)
715 331–335.
- 716 [42] G.M. Neumann, P.G. Cullis, P.J. Derrick, Mass Spectrometry of Polymers:
717 Polypropylene Glycol, *Zeitschrift Für Naturforschung A.* 35 (1980).
718 doi:10.1515/zna-1980-1015.
- 719 [43] K. Yang, X. Han, Accurate quantification of lipid species by electrospray
720 ionization mass spectrometry - Meet a key challenge in lipidomics, *Metabolites.* 1
721 (2011) 21–40.

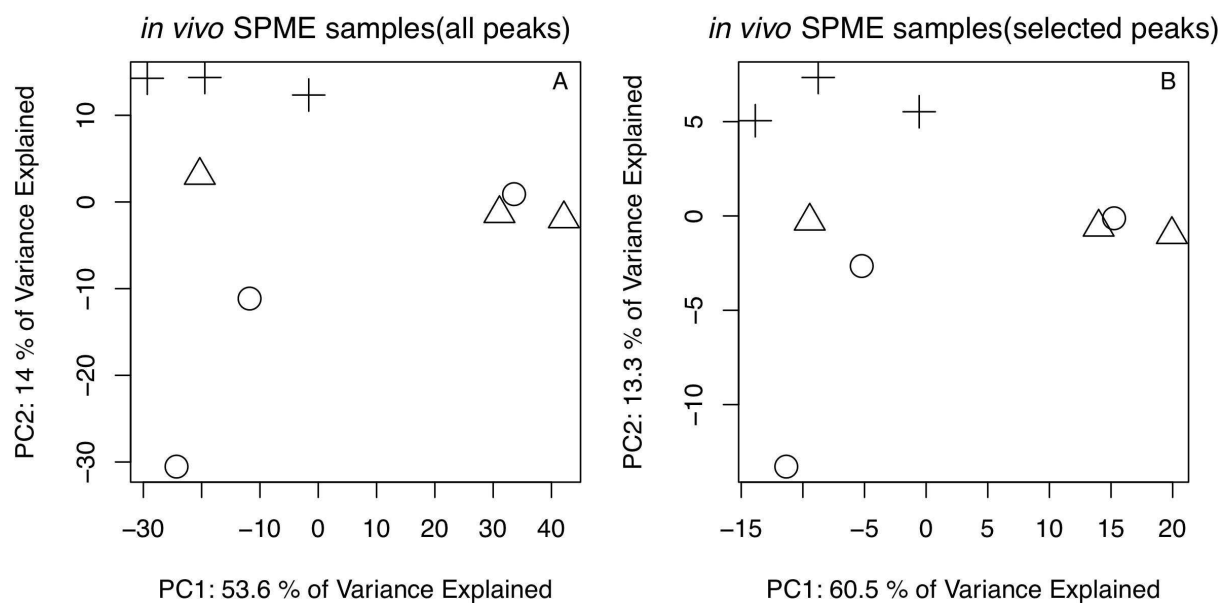


722
 723
 724
 725
 726
 727
 728
 729
 730
 731
 732

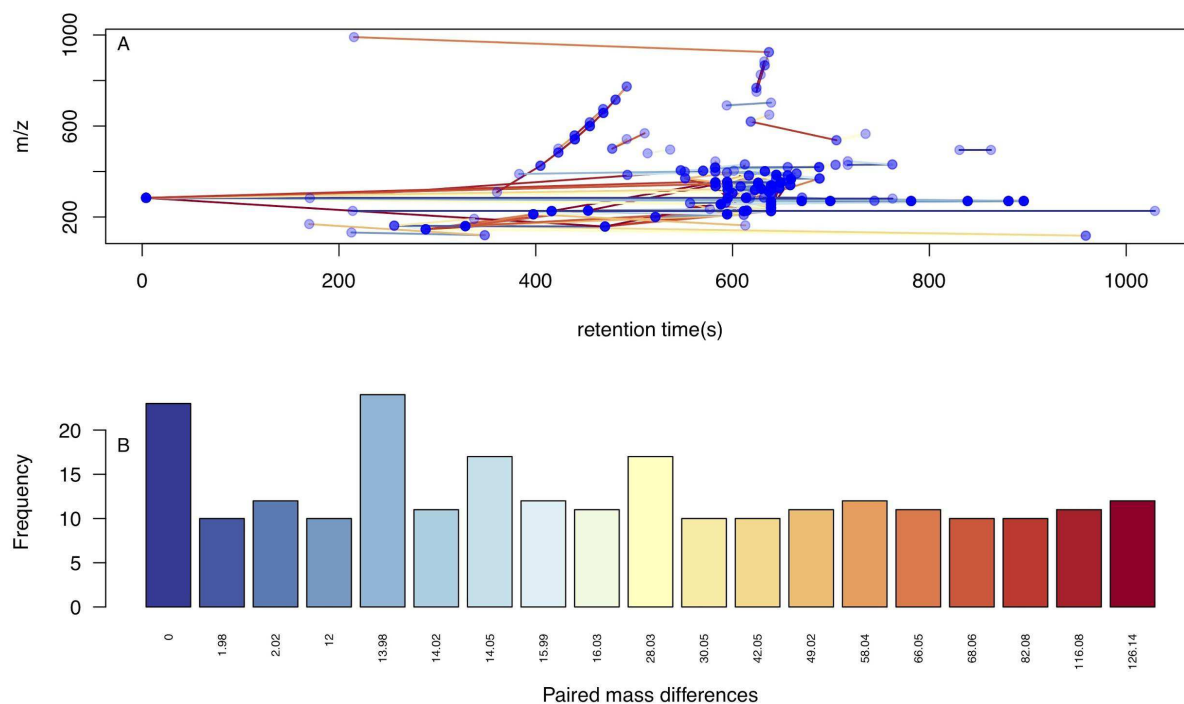
Scheme 1. Demonstration of GlobalStd algorithm. Different colors stand for peaks from different retention time groups. Blue stands for mass pairs with high frequency paired mass distances (PMD). Step 1 indicates the retention time cluster analysis to find Pseudospectra for potential compounds. Step 2 indicates the PMD-based global search. Step 3 indicates selection of independent peaks and detailed process is shown on the right flowchart.



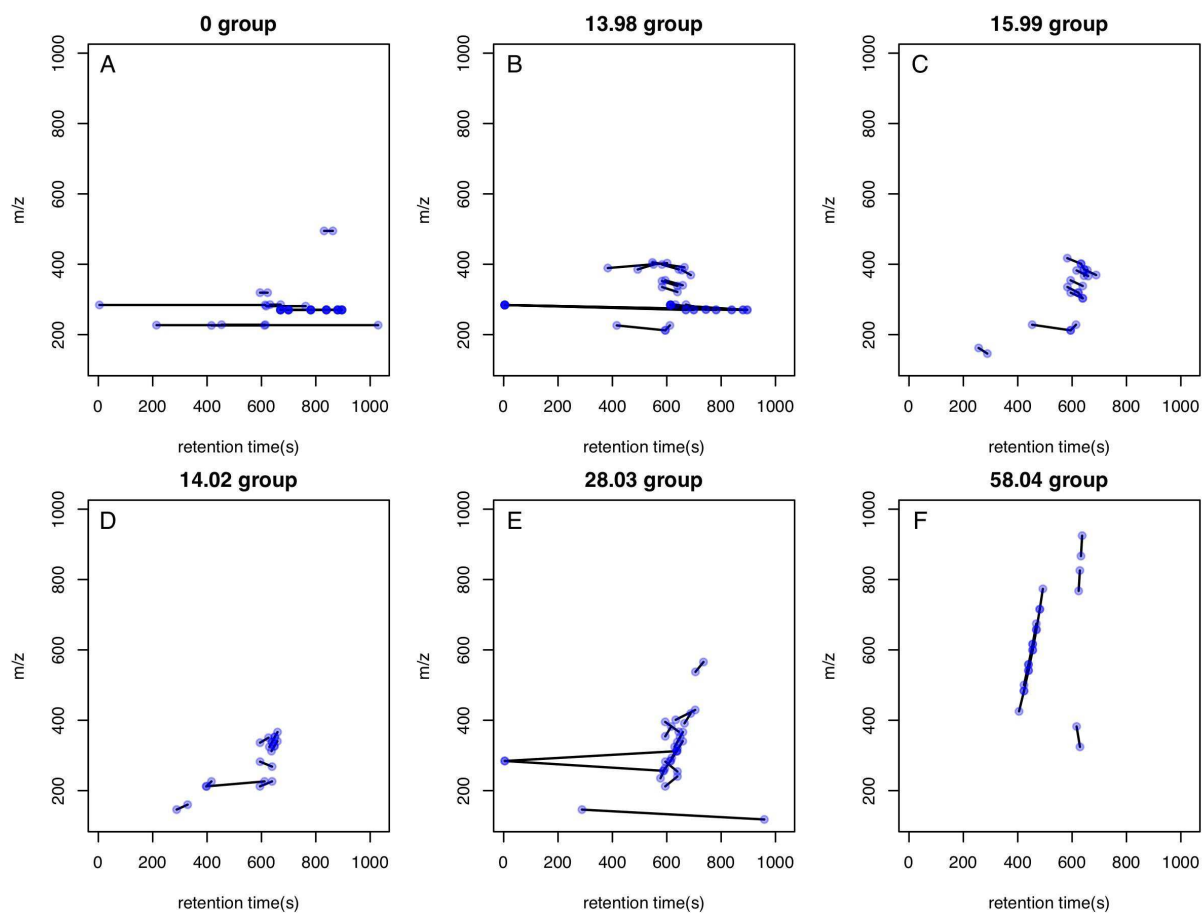
733
 734 Figure 1. Paired mass distance (PMD) analysis for *in vivo* SPME sampling data. The plot
 735 on the top of the figure(A) illustrates the PMD relationship across retention time – m/z
 736 profile, while the bottom plot(B) shows the corresponding PMD frequency. The paired
 737 relationships are reflected by the lines between peaks. The colors of the segments in the
 738 top plot correspond to the colors in the bar plot below, indicating the PMD groups.
 739
 740
 741



742
743 Figure 2. Score plot from principal components analysis (PCA) of raw peaks (A) and Std
744 peaks selected by GlobalStd algorithm (B). Use of the same symbols in plots indicates
745 technique replicates, while different symbols indicate biological replicates.
746



747
748 Figure 3. Structure/reaction directed analysis in *in vivo* SPME sampling data for peaks
749 selected via application of the GlobalStd algorithm. The plot above (A) illustrates the
750 PMD relationship across retention time and m/z profile, while the bottom plot (B)
751 indicates the corresponding PMD frequency. The paired relationships are reflected by the
752 line between peaks. The colors of the segments in figure 3A correspond the colors in the
753 bar plot below, which categorize the PMD groups.
754
755



756
757 Figure 4. m/z - retention time peak profiles obtained via GlobalStd algorithm analysis for
758 PMD 0Da (A, isomers), PMD 13.98Da (B, replacement of oxygen atom and two
759 hydrogen atoms), PMD 15.99Da (C, oxidation), PMD 14.02Da (D, homologous series
760 with $-\text{CH}_2-$), PMD 28.03Da (E, homologous series with $-\text{C}_2\text{H}_4-$), and PMD 58.04Da (F,
761 homologous series with $-\text{C}_3\text{H}_6\text{O}-$). The paired relationships are reflected by the lines
762 between peaks.

763
764
765

Highlights

- Algorithms were developed to reduce redundant peaks in metabolomics data profile
- 20% of the original peaks could stand for the major variances for data
- Quantitative analysis could be performed at structure/reaction level
- Unknown structure/reaction relationships could be revealed by *in vivo* SPME sampling