

1 [TITLE] Incorporating in-source fragment information improves metabolite identification  
2 accuracy in untargeted LC-MS datasets

3 *[AUTHOR NAMES] Phillip M. Seitzer<sup>1</sup>, Brian C. Searle<sup>1,2,\*</sup>*

4 [AUTHOR ADDRESS]

5 <sup>1</sup>Proteome Software, 1340 SW Bertha Blvd Ste 10, Portland, OR

6 <sup>2</sup>Department of Genome Sciences, University of Washington, Seattle, WA

7 \*Corresponding author, email: [brian.searle@proteomesoftware.com](mailto:brian.searle@proteomesoftware.com)

8 [KEYWORDS] Metabolomics, In-source fragment, Identification, Scoring, MS1  
9 Spectrum, Untargeted, Spectral Library

10 [ABSTRACT]

11 In-source fragmentation occurs as a byproduct of electrospray ionization. We find that  
12 ions produced as a result of in-source fragmentation often match fragment ions produced  
13 during MS/MS fragmentation and we take advantage of this phenomenon in a novel  
14 algorithm to analyze LC-MS metabolomics datasets. Our approach organizes co-eluting  
15 MS1 features into a single peak group and then identifies in-source fragments among co-  
16 eluting features using MS/MS spectral libraries. We tested our approach using  
17 previously published data of verified metabolites, and compared the results to features  
18 detected by other mainstream metabolomics tools. Our results indicate that considering  
19 in-source fragment information as a part of the identification process increases annotation

20 quality, allowing us to leverage MS/MS data in spectrum libraries even if MS/MS scans  
21 were not collected.

22 [TEXT]

## 23 **INTRODUCTION**

24 Confidently identifying metabolites in LC-MS metabolomics datasets is a  
25 challenging problem<sup>1</sup>. Both targeted<sup>2-3</sup> and untargeted<sup>4</sup> LC-MS raw data can be  
26 internally or externally calibrated with chemical standards. While this can yield highly  
27 accurate metabolite detections, the approach is constrained to only measure endogenous  
28 levels of those standard metabolites. Additionally, external calibrant data must be  
29 reacquired when chromatographic conditions or instrument settings change, making it  
30 potentially prohibitively expensive and time-consuming to produce.

31 When internal or external standards are unavailable, metabolomics studies  
32 typically leverage several independent lines of evidence to detect metabolites, including  
33 accurate mass, retention time, and agreement between observed and theoretical isotopic  
34 peak intensities. Different types of identification information may be aggregated to  
35 produce a single identification score<sup>5</sup>, or identification probabilities using Bayesian  
36 networks<sup>6</sup> and target-decoy approaches<sup>7</sup>. A popular alternative for analyzing untargeted  
37 LC-MS/MS data is matching acquired MS/MS against one or more large spectral  
38 libraries, such as NIST<sup>8</sup>, HMDB<sup>9</sup>, and METLIN<sup>10</sup>. While the number of features without  
39 MS/MS spectra acquired using data dependent acquisition (DDA) experiments remains  
40 significant, efforts to increase the number of MS1 features fragmented by the mass  
41 spectrometer<sup>11</sup> and applications of data independent acquisition<sup>12-13</sup> may improve data  
42 consistency.

43           However, many metabolomics experiments are still collected using LC-MS only,  
44 and even in LC-MS/MS datasets, many features only contain MS1 information. Without  
45 MS/MS information, search engines can only use accurate mass and isotopic distributions  
46 based on molecular formulae to detect metabolites<sup>14</sup>. As many metabolites share  
47 molecular formulae, scanning MS1-only data against spectral libraries yields incomplete,  
48 ambiguous, or partial metabolite identifications. Additionally, when individual  
49 metabolites ionize, they can produce unanticipated MS1 features as a result of neutral  
50 losses, in-source fragmentation, multimerization, and adducts<sup>12,15</sup>, further complicating  
51 the annotation process.

52           Here we present an approach to identify metabolites in untargeted LC-MS data by  
53 identifying in-source fragments that match to fragment peaks in MS/MS spectral  
54 libraries. To accomplish this, we have developed an algorithm to form consensus MS1  
55 peak groups from a set of raw data files and use those peak groups in library searching.  
56 We have tested our method by comparing the feature detection, deisotoping and grouping  
57 steps of our algorithm to two mainstream open-source approaches using a complex LC-  
58 MS dataset containing 75 verified compounds. We find that our feature detection,  
59 deisotoping and peak grouping steps identify more of the verified compound features  
60 than other approaches. We also find that identifying in-source fragments in LC-MS data  
61 and including this information as a part of our identification process improves the  
62 accuracy of metabolite identifications.

63

64   **EXPERIMENTAL PROCEDURES**

65 We downloaded mzML raw data from the Metabolights study 67 (MTBLS67)<sup>16</sup>  
66 from the Metabolights raw data portal<sup>17</sup>. We processed raw files with MSConvertGUI  
67 (Proteowizard version 3.0.9987)<sup>18</sup> to strip them of MS/MS scans, and generated both a  
68 centroided set and an uncentroided set of sample files (using the parameters cwt  
69 centroiding, snr = 0.1, peakSpace = 0.1).

70 We independently processed the uncentroided positive and negative mode files  
71 using Scaffold Elements 2.0.0 with search parameters that were chosen to match the  
72 original MTBLS67 study (specific search parameters are listed in **Supporting**  
73 **Information Table 1**). Monoisotopic peaks were searched against the NIST 2017<sup>8</sup> and  
74 METLIN<sup>10</sup> spectral libraries, as well as an empty library (to generate a baseline list of all  
75 detected features). We also generated an R script (**Supporting Information Script 1**)  
76 using Bioconductor<sup>19</sup> to drive XCMS<sup>20-21</sup> (version 3.0.2) and CAMERA<sup>22</sup> (version  
77 1.34.0). The script performed peak detection, peak grouping, and isotope detection on  
78 both the uncentroided sample files (using XCMS “matchedFilter”<sup>20</sup>) and the centroided  
79 sample files (using XCMS “centWave”<sup>21</sup>). We analyzed positive mode and negative  
80 mode files separately using search parameters that were chosen to match the original  
81 study (specific search parameters are listed in **Supporting Information Script 1**). The  
82 *m/z* and retention time coordinates of the 75 verified metabolites were compared to all  
83 monoisotopic *m/z* and retention time features identified by XCMS-matchedFilter +  
84 CAMERA, XCMS-centWave + CAMERA, and Scaffold Elements (script available in  
85 **Supporting Information Script 2**).

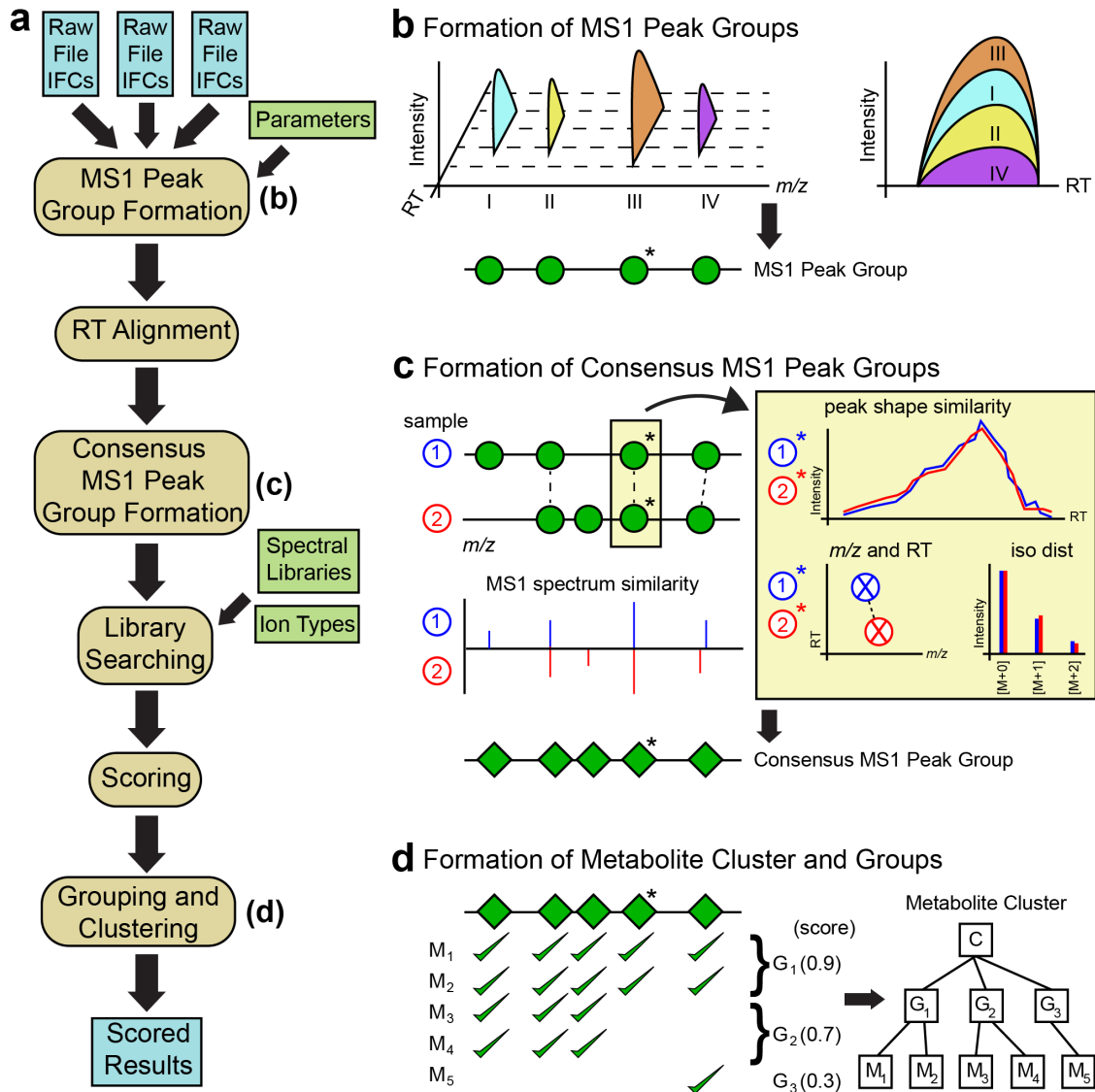
86

87 **RESULTS AND DISCUSSION**

## 88 **The Scaffold Elements algorithmic workflow**

89 We have developed an automated workflow to identify metabolites from  
90 untargeted LC-MS raw files using spectral libraries (**Fig 1a**). Briefly (see **Supporting**  
91 **Information Note 1** and **Supporting Information Figure 1** for further details), we first  
92 organize raw data into isotopic feature clusters (IFCs) that contain a monoisotopic [M+0]  
93 feature, and [M+1] and [M+2] isotopic features. IFCs from the same sample are formed  
94 into MS1 peak groups based on elution profile (**Fig 1b**). This step ensures that all ions  
95 produced during ionization of a single metabolite remain organized together. Failure to  
96 properly account for ionization effects can lead to ion misannotation, especially of in-  
97 source fragments<sup>23</sup>. We then align MS1 peak groups from all samples in the experiment  
98 to form cross-sample consensus MS1 peak groups. The formation of consensus elements  
99 is based on a number of independent metrics, including MS1 spectral similarity, peak  
100 shape, and agreement in  $m/z$  and retention time (**Fig 1c**). Finally, we search consensus  
101 MS1 peak groups against spectral libraries and score metabolite groups and clusters (**Fig**  
102 **1d**). Score values increase both with agreement (higher mass accuracy and agreement  
103 with theoretically predicted isotopic distributions) and the amount of evidence associated  
104 with a metabolite annotation (number of ion types and in-source fragments identified).

105 **Figure 1:**



106

107

## 108 Development of a “gold-standard” MS1-only dataset

109 We benchmarked our approach using the Metabolights study 67 (MTBLS67)<sup>16</sup>.

110 This study identified and quantified 75 yeast metabolites from nitrogen-starved

111 *Saccharomyces Pombe* whole cell lysates using DDA-based LC-MS/MS. Sajiki et al

112 confirmed the MS/MS fragmentation patterns and retention times of these metabolites

113 using external standards. In an effort to produce a “gold-standard” MS1-only dataset of a

114 complex metabolome with endogenous targets, we stripped these raw files of MS/MS  
115 scans. This produced a mock MS1-only data set containing 75 independently verified  
116 compounds.

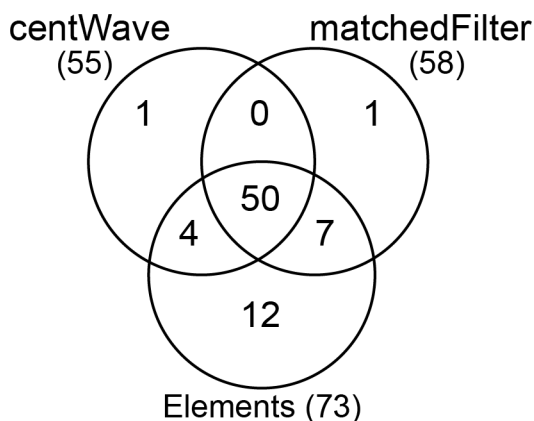
117

### 118 **Comparing peak detection algorithms**

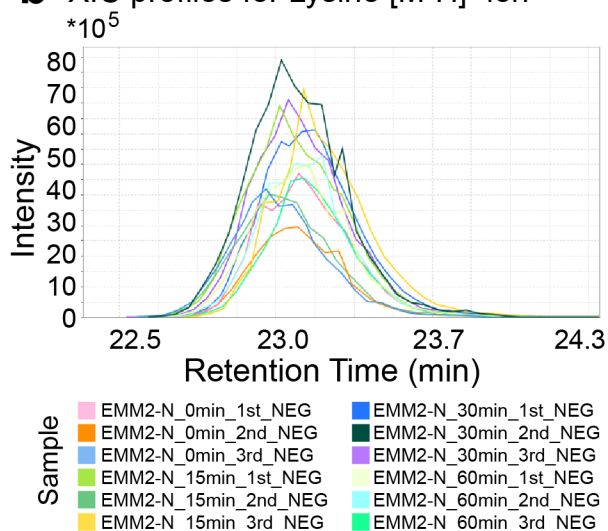
119 We compared the peak detection, isotopic clustering, and peak grouping steps of  
120 our approach to two XCMS-based workflows, either XCMS “matchedFilter”<sup>20</sup> or XCMS  
121 “centWave”<sup>21</sup> peak detection, both followed by CAMERA isotopic grouping<sup>22</sup>. Scaffold  
122 Elements was executed without library matching to generate a list of all dataset  
123 features. We found that Scaffold Elements was able to detect more of the features  
124 associated with verified metabolites than either XCMS-CAMERA workflow, including  
125 12 that were not identified by either approach (**Fig 2a**). However, since Scaffold  
126 Elements reported more features than either XCMS-CAMERA workflow (**Supporting**  
127 **Information Figure 2**), we were concerned that there would be a higher chance of noise  
128 matching a verified metabolite *m/z* and retention time coordinate by chance. To ensure  
129 that Scaffold Elements returned well-formed peaks, we manually investigated the  
130 features associated with the 12 metabolites that were only identified by Scaffold  
131 Elements. We found that 11 of these 12 verified metabolite features had a clear,  
132 reproducible signal (**Supporting Information Figure 3**). Extracted ion chromatograms  
133 of features corresponding to one representative verified metabolite (Lysine) are shown in  
134 **Fig 2b**.

135 **Figure 2:**

**a** Verified Metabolite Feature Recall



**b** XIC profiles for Lysine [M-H]<sup>-</sup> ion



136

137

138 **Using in-source fragments in scoring improves annotation quality**

139 We next aimed to determine if searching for in-source fragments in MS1 peak  
140 groups improved metabolite annotation quality. We searched the MTBLS67 sample files  
141 with the NIST<sup>8</sup> and METLIN<sup>10</sup> spectral libraries, which together contained 65 of the 75  
142 verified metabolites (**Supporting Information Table 2**). Our feature detection  
143 algorithm identified the correct *m/z* and retention time feature for 63 of these 65  
144 metabolites. However, multiple library annotations were returned for these features.



145 Scaffold Elements' scoring algorithm organized these annotations into clusters of  
146 metabolite groups, and ranked the annotations within each metabolite group.

147 We evaluated metabolite detection performance based on three metrics. For each  
148 independent search, we determined the proportion of correct annotations (where the  
149 annotation had the highest score in the metabolite group), unambiguous annotations  
150 (where the correct annotation had a uniquely higher score than all other annotations in the  
151 metabolite group), and unmistakable annotations (where the correct annotation was the  
152 only annotation in the metabolite group). Our approach of incorporating in-source  
153 fragment information in scoring improved all three of these metrics, notably increasing  
154 the proportion of unambiguous and unmistakable annotations by 22% and 60%,  
155 respectively (**Table 1**). In many cases, the inclusion of in-source fragments in the search  
156 yielded rich MS1 peak groups that matched multiple MS/MS fragment peaks from the  
157 corresponding library spectrum with high mass accuracy (**Fig 3**).

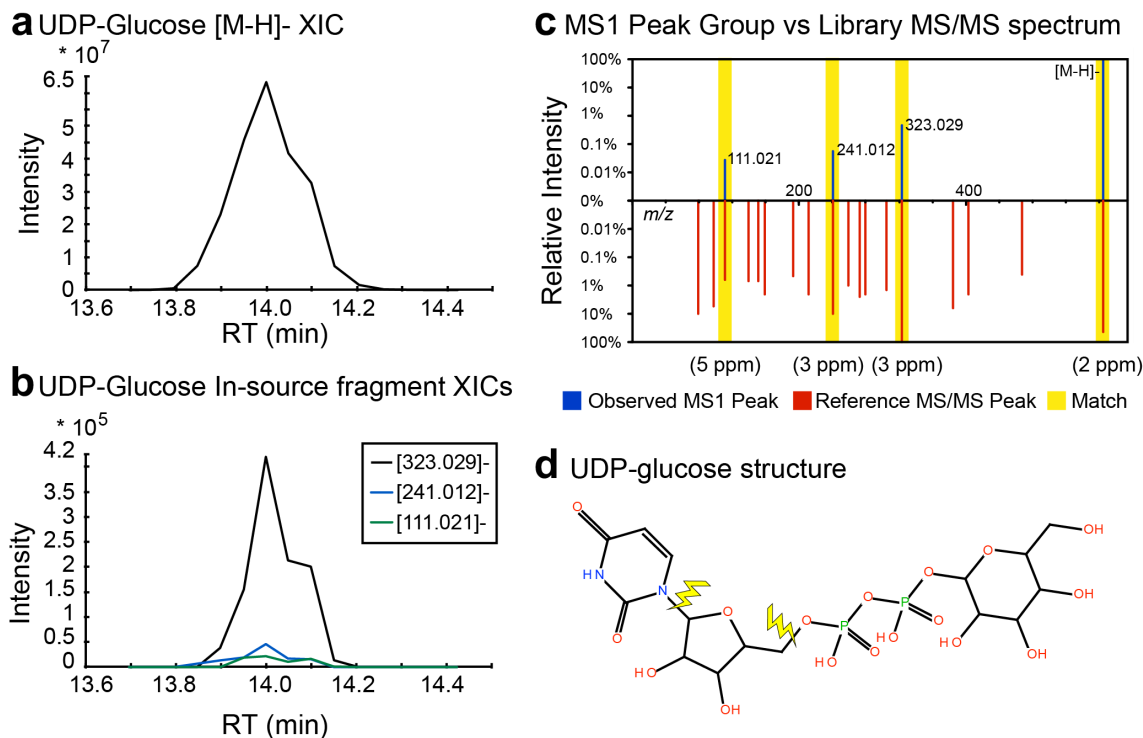
158

159 **Table 1:**

Search	Correct <sup>a</sup>	Unambiguous <sup>b</sup>	Unmistakable <sup>c</sup>
<b>Major Ion</b>	96.8%	59.7%	40.3%
<b>Major Ion + ISFs</b>	98.4%	72.6%	64.5%

160

161 **Figure 3:**



162

## 163 CONCLUSIONS

164 We have developed an approach to account for ionization effects by forming  
165 consensus MS1 peak groups prior to spectral library matching, and to use in-source  
166 fragments in those groups to perform pseudo-MS/MS library searching. Our results  
167 indicate that considering in-source fragments as part of the identification process  
168 improves confidence in metabolite detections. To increase the availability of these  
169 algorithms, we have made this tool available as a module in the Scaffold Elements  
170 software package distributed by Proteome Software.

171 Our results also demonstrate a caveat of spectral library search-based approaches:  
172 it is only possible to identify metabolites that are present in the specific spectral library  
173 (or libraries) searched. In our case, only 65 (86.7%) of the verified compounds were  
174 present in the NIST and METLIN spectral libraries (**Supporting Information Table**  
175 **2**). If a compound is present in the data but absent from the library, the compound will

176 either be misidentified or remain unidentified. Without prior knowledge of which  
177 compounds are actually contained in the data, we can use our scoring approach to  
178 determine which annotations correspond to real compounds and which are  
179 misidentifications. We believe that improving candidate scoring is particularly important  
180 for analyzing untargeted metabolomics LC-MS data, as the ground truth identification  
181 might be absent from the library.

182

183 [FIGURES]

184 **Figure 1. Scaffold Elements metabolite identification and scoring algorithm**

185 **(a)** Complete workflow of Scaffold Elements identification and scoring algorithm. Tan,  
186 rounded boxes indicate algorithmic steps, green boxes indicate user-specified inputs, and  
187 blue boxes indicate algorithmic inputs and outputs. **(b)** An MS1 peak group is formed in  
188 a single sample by combining four co-eluting isotopic feature clusters (IFCs) (I, II, III,  
189 and IV). IFCs are represented as green circles on a line, with an asterisk indicating the  
190 most intense IFC in the peak group. **(c)** A consensus MS1 peak group is formed by  
191 comparing MS1 peak groups from each sample. A cross-sample MS1 spectrum  
192 similarity score is evaluated considering all IFCs in each peak group, and additional  
193 comparisons are made between a representative IFC from each MS1 peak group  
194 individually (light yellow boxes). The resulting consensus MS1 peak group is represented  
195 as green diamonds on a line, with an asterisk indicating the most intense consensus IFC  
196 in the consensus MS1 peak group. **(d)** Multiple putatively identified metabolites are  
197 organized into groups and clusters based on the consensus IFCs within a consensus MS1  
198 peak group. In this schematic, a consensus MS1 spectrum of five IFCs was identified by

199 five metabolites, which were organized into a cluster containing three groups, one of  
200 which contained only a single metabolite. Identification scores (shown next to each  
201 group in parentheses) indicate the most likely metabolite annotation for this cluster.

202

### 203 **Figure 2: Scaffold Elements feature detection comparison**

204 **(a)** Comparison of verified metabolite features identified by XCMS-centWave +  
205 CAMERA (centWave), XCMS-matchedFilter + CAMERA (matchedFilter) and Scaffold  
206 Elements (elements). Scaffold Elements identified 73 of the 75 features associated with  
207 verified metabolites, including 12 that were not detected by either XCMS-CAMERA  
208 workflow. **(b)** Extracted ion chromatograms (XICs) of a verified metabolite ion for  
209 Lysine ([M-H]<sup>-</sup> ion), which was identified only in Elements. The overlay plot of XICs  
210 shows a reasonable peak shape for this ion, which was independently identified in all 12  
211 negative mode samples and correctly organized together into a single feature group.

212

### 213 **Figure 3: UDP-glucose MS1 peak group**

214 **(a)** An extracted ion chromatogram (XIC) of [M-H]<sup>-</sup> ion of UDP-glucose and **(b)** XICs of  
215 three detected in-source fragment peaks. **(c)** A butterfly plot comparing observed MS1  
216 peak group of UDP-glucose ([M-H]<sup>-</sup> ion and three in-source fragment peaks) to METLIN  
217 library spectrum ID:6698 (METLIN ID 3598). Intensities are shown as a relative  
218 percentage to max spectral peak on a logarithmic scale to allow visualization of low-  
219 intensity peaks. The mass tolerance in ppm for each peak match is shown below butterfly  
220 plot. **(d)** The structure of UDP-glucose, with fragmentation sites corresponding to

221 observed in-source fragments indicated by yellow lightning bolts. All observed data in  
222 figure was taken from the sample “EMM2-N\_0min\_2nd\_NEG”.

223

224 [TABLES]

225 Table 1. Annotation of verified metabolites with and without consideration of in-source  
226 fragmentation (ISF) events in the identification process. <sup>a</sup>The annotation had the highest  
227 score. <sup>b</sup>The correct annotation had a uniquely higher score than all other annotations.  
228 <sup>c</sup>The correct annotation was the only annotation in the metabolite group.

229

230 [ASSOCIATED CONTENT]

231 **Supporting Information Note 1: Detailed description of Scaffold Elements 2.0**  
232 **metabolite identification and scoring algorithm.** A detailed description of the Scaffold  
233 Elements 2.0 metabolite identification and scoring algorithm. Also includes a description  
234 of feature finding and isotopic grouping.

235 **Supporting Information Figure 1: Feature finding algorithm** Diagram of major steps  
236 of Scaffold Elements feature detection algorithm.

237 **Supporting Information Figure 2: Number of features identified by different**  
238 **programs** Summary of number of features identified by XCMS-centWave + CAMERA,  
239 XCMS-matchedFilter + CAMERA, and Scaffold Elements.

240 **Supporting Information Figure 3: XICs of verified features only detected by**  
241 **Scaffold Elements** Description and summary of 12 verified features only detected by  
242 Scaffold Elements, including overlaid XICs (showing XIC of each feature in all samples  
243 where it was detected).

244 **Supporting Information Table 1: Scaffold Elements parameters** Table of parameter  
245 used in all Scaffold Elements analyses.

246 **Supporting Information Table 2: Detailed Results of in-source fragment annotation**  
247 **comparison analysis** Detailed summary of the annotation results for 75 verified  
248 metabolites with and without consideration of in-source fragments.

249 **Supporting Information Script 1: XCMS CAMERA workflows (R script).** R Script  
250 for generating (*m/z*, RT) feature list files using both XCMS-matchedFilter on profile  
251 mode files and XCMS-centWave on centroided files. Uses Bioconductor, XCMS, and  
252 CAMERA (for isotopic grouping).

253 **Supporting Information Script 2: Comparison of XCMS CAMERA workflows vs**  
254 **Scaffold Elements (Java script)** Java script comparing output of Scaffold Elements and  
255 XCMS-CAMERA workflows to features corresponding to verified metabolites.

256 [AUTHOR INFORMATION]

257 **Corresponding Author**

258 Brian C. Searle, [brian.searle@proteomesoftware.com](mailto:brian.searle@proteomesoftware.com)

259 **Author Contributions**

260 The study was conceived by B.C.S. and P.M.S. The algorithm was implemented and  
261 evaluated by P.M.S. P.M.S. and B.C.S. wrote the paper. All authors have given approval  
262 to the final version of the manuscript.

263 [ACKNOWLEDGEMENT]

264 We would like to acknowledge the entire staff at Proteome Software, Inc. for fruitful  
265 scientific discussions and feedback associated with development and implementation of  
266 the algorithm.

267 [ABBREVIATIONS]

268 LC-MS, liquid chromatography mass spectrometry, LC-MS/MS, liquid chromatography  
269 tandem mass spectrometry, MS1, mass spectrometry, MS/MS, tandem mass  
270 spectrometry, IFC isotopic feature cluster, NIST, national institute of standards and  
271 technology, HMDB, human metabolome database, MTBLS, Metabolights, DDA, data-  
272 dependent acquisition, DIA, data-independent acquisition, RT, retention time. ISF, in-  
273 source fragment, IFC, isotopic feature cluster, XIC, extracted ion chromatogram

274 [REFERENCES]

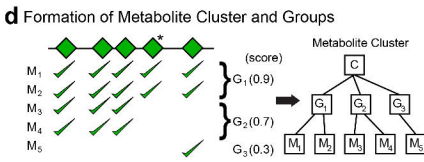
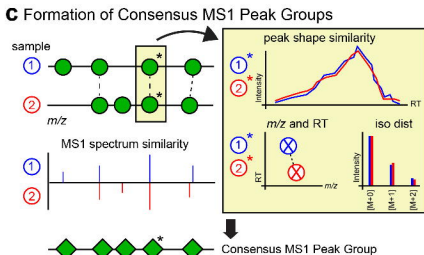
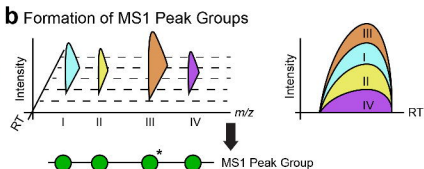
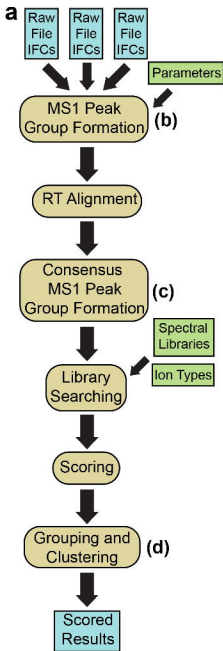
- 275 (1) Schrimpe-Rutledge, Alexandra C., et al. "Untargeted metabolomics strategies—  
276 challenges and emerging directions." *Journal of The American Society for Mass*  
277 *Spectrometry* **2016**, 27(12), 1897-1905.
- 278 (2) Šimura, J., Antoniadi, I., Široká, J., Tarkowská, D., Strnad, M., Ljung, K., &  
279 Novák, O. (2018). Plant hormonomics: multiple phytohormone profiling by  
280 targeted metabolomics. *Plant physiology*, 177(2) **2018**, 476-489.
- 281 (3) Jacob, M., Malkawi, A., Albast, N., Al Bougha, S., Lopata, A., Dasouki, M., &  
282 Rahman, A. M. A. (2018). A targeted metabolomics approach for clinical  
283 diagnosis of inborn errors of metabolism. *Analytica chimica acta* **2018**, 1025,  
284 141-153.

- 285 (4) Patti, Gary J., Oscar Yanes, and Gary Siuzdak. "Innovation: Metabolomics: the  
286 apogee of the omics trilogy." *Nature reviews Molecular cell biology* **2012**, 13(4),  
287 13.4, 263.
- 288 (5) Creek, D. J.; Dunn, W. B.; Fiehn, O.; Griffin, J. L.; Hall, R. D.; Lei, Z.; Mistrik,  
289 R.; Neumann, S.; Schymanski, E. L.; Sumner, L. W.; et al. Metabolite  
290 identification: are you sure? And how do your peers gauge your confidence?  
291 *Metabolomics* **2014**, 10, 350–353.
- 292 (6) Daly, R.; Rogers, S.; Wandy, J.; Jankevics, A.; Burgess, K. E. V; Breitling, R.  
293 MetAssign: Probabilistic annotation of metabolites from LC-MS data using a  
294 Bayesian clustering approach. *Bioinformatics* **2014**, 30 (19), 2764–2771.
- 295 (7) Wang, X., Jones, D. R., Shaw, T. I., Cho, J. H., Wang, Y., Tan, H., ... & Peng, J.  
296 (2018). Target-decoy Based False Discovery Rate Estimation for Large-scale  
297 Metabolite Identification. *Journal of proteome research*. **2018**, DOI:  
298 10.1021/acs.jproteome.8b00019
- 299 (8) Yang, X., Neta, P., & Stein, S. E. (2017). Extending a Tandem Mass Spectral  
300 Library to Include MS2 Spectra of Fragment Ions Produced In-Source and MSn  
301 Spectra. *Journal of The American Society for Mass Spectrometry* **2017**, 28(11),  
302 2280-2287.
- 303 (9) Wishart, D. S., Feunang, Y. D., Marcu, A., Guo, A. C., Liang, K., Vázquez-  
304 Fresno, R., ... & Sayeeda, Z. HMDB 4.0: the human metabolome database for  
305 2018. *Nucleic acids research* **2017**, 46(D1), D608-D617.

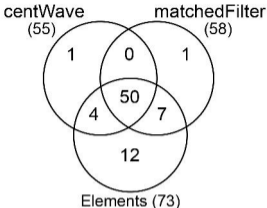


- 306 (10) Guijas, C., Montenegro-Burke, J. R., Domingo-Almenara, X., Palermo, A.,  
307 Warth, B., Hermann, G., ... & Wolan, D. W. METLIN: a technology platform for  
308 identifying knowns and unknowns. *Analytical chemistry*, **2018** 90(5), 3156-3164.
- 309 (11) Broeckling, C. D., Hoyes, E., Richardson, K., Brown, J. M., & Prenni, J. E.  
310 DataSet-Dependent Acquisition enables comprehensive tandem mass  
311 spectrometry coverage of complex samples. *Analytical chemistry*. **2018** DOI:  
312 10.1021/acs.analchem.8b00929
- 313 (12) Broeckling, C. D., Afsar, F. A., Neumann, S., Ben-Hur, A., & Prenni, J. E.  
314 RAMClust: a novel feature clustering method enables spectral-matching-based  
315 annotation for metabolomics data. *Analytical chemistry* **2014**, 86(14), 6812-6817.
- 316 (13) Zhou, J., Li, Y., Chen, X., Zhong, L., & Yin, Y. Development of data-  
317 independent acquisition workflows for metabolomic analysis on a quadrupole-  
318 orbitrap platform. *Talanta* **2017**, 164, 128-136
- 319 (14) Kind, Tobias, and Oliver Fiehn. "Metabolomic database annotations via query of  
320 elemental compositions: mass accuracy is insufficient even at less than 1 ppm."  
321 *BMC bioinformatics* **2006** 7(1), 234.
- 322 (15) Keller, B. O., Sui, J., Young, A. B., & Whittall, R. M. Interferences and  
323 contaminants encountered in modern mass spectrometry. *Analytica chimica acta*  
324 **2008**, 627(1), 71-81.
- 325 (16) Sajiki, K.; Pluskal, T.; Shimanuki, M.; Yanagida, M. Metabolomic Analysis of  
326 Fission Yeast at the Onset of Nitrogen Starvation. *Metabolites* **2013**, 3, 1118-  
327 1129.

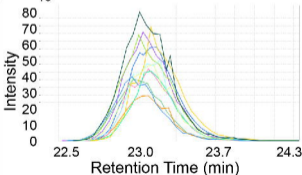
- 328 (17) Haug, K.; Salek, R. M.; Conesa, P.; Hastings, J.; De Matos, P.; Rijnbeek, M.;  
329 Mahendraker, T.; Williams, M.; Neumann, S.; Rocca-Serra, P.; et al.  
330 MetaboLights - An open-access general-purpose repository for metabolomics  
331 studies and associated meta-data. *Nucleic Acids Res.* **2013**, 41 (11), 1–6.
- 332 (18) Kessner, Darren, et al. "ProteoWizard: open source software for rapid  
333 proteomics tools development." *Bioinformatics* **2008**, 24(21), 2534-2536.
- 334 (19) Huber, Wolfgang, et al. "Orchestrating high-throughput genomic analysis with  
335 Bioconductor." *Nature methods* **2015**, 12(2), 115.
- 336 (20) Smith, Colin A., et al. "XCMS: processing mass spectrometry data for  
337 metabolite profiling using nonlinear peak alignment, matching, and  
338 identification." *Analytical chemistry* **2006**, 78(3), 779-787.
- 339 (21) Tautenhahn, Ralf, Christoph Boettcher, and Steffen Neumann. "Highly sensitive  
340 feature detection for high resolution LC/MS." *BMC bioinformatics* **2008**, 9(1),  
341 504.
- 342 (22) Kuhl, Carsten, et al. "CAMERA: an integrated strategy for compound spectra  
343 extraction and annotation of liquid chromatography/mass spectrometry data sets."  
344 *Analytical chemistry* **2011**, 84(1), 283-289.
- 345 (23) Xu, Y. F., Lu, W., & Rabinowitz, J. D. Avoiding misannotation of in-source  
346 fragmentation products as cellular metabolites in liquid chromatography–mass  
347 spectrometry-based metabolomics. *Analytical chemistry* **2015**, 87(4), 2273-2281.

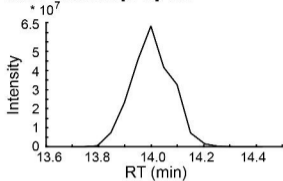
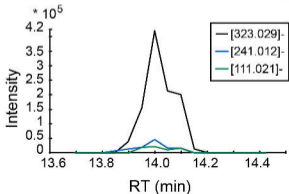
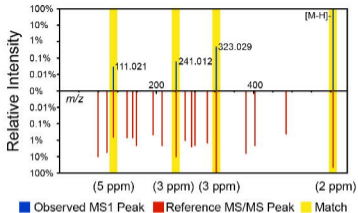


## a Verified Metabolite Feature Recall



## b XIC profiles for Lysine [M-H]<sup>-</sup> ion



**a** UDP-Glucose [M-H]<sup>-</sup> XIC**b** UDP-Glucose In-source fragment XICs**c** MS1 Peak Group vs Library MS/MS spectrum**d** UDP-glucose structure