

# RAMClust: A Novel Feature Clustering Method Enables Spectral-Matching-Based Annotation for Metabolomics Data

C. D. Broeckling,<sup>\*,†</sup> F. A. Afsar,<sup>\*,‡</sup> S. Neumann,<sup>\*,∇</sup> A. Ben-Hur,<sup>\*,‡</sup> and J. E. Prenni<sup>\*,†,§</sup>

<sup>†</sup>Proteomics and Metabolomics Facility, Colorado State University, Fort Collins, Colorado 80523, United States

<sup>‡</sup>Department of Computer Science, Colorado State University, Fort Collins, Colorado 80523, United States

<sup>§</sup>Department of Biochemistry, Colorado State University, Fort Collins, Colorado 80523, United States

<sup>∇</sup>Department of Stress and Developmental Biology, Leibniz Institute of Plant Biochemistry, 06108 Halle, Germany

**ABSTRACT:** Metabolomic data are frequently acquired using chromatographically coupled mass spectrometry (MS) platforms. For such datasets, the first step in data analysis relies on feature detection, where a feature is defined by a mass and retention time. While a feature typically is derived from a single compound, a spectrum of mass signals is more a more-accurate representation of the mass spectrometric signal for a given metabolite. Here, we report a novel feature grouping method that operates in an unsupervised manner to group signals from MS data into spectra without relying on predictability of the in-source phenomenon. We additionally address a fundamental bottleneck in metabolomics, annotation of MS level signals, by incorporating indiscriminant MS/MS (idMS/MS) data implicitly: feature detection is performed on both MS and idMS/MS data, and feature–feature relationships are determined simultaneously from the MS and idMS/MS data. This approach facilitates identification of metabolites using in-source MS and/or idMS/MS spectra from a single experiment, reduces quantitative analytical variation, compared to single-feature measures, and decreases false positive annotations of unpredictable phenomenon as novel compounds. This tool is released as a freely available R package, called RAMClustR, and is sufficiently versatile to group features from any chromatographic-spectrometric platform or feature-finding software.



Mass spectrometry (MS) has long been utilized for detecting and quantifying small molecules, particularly when coupled to separation tools such as gas chromatography (GC), liquid chromatography (LC), or capillary electrophoresis (CE). The strengths of these chromatographically coupled mass spectrometry platforms have been leveraged toward global metabolite profiling approaches, or metabolomics. The development of electrospray ionization (ESI)<sup>1</sup> was an important technological milestone, which allowed for the coupling of liquid separation methods to mass spectrometers. This development obviated the volatility requirement imposed by gas chromatography and supported development and expansion of both metabolomics and proteomics. Electrospray is considered a “soft” ionization technique, by which the molecular ion of the compound is generally more dominant than that achieved using “hard” ionization methods such as electron impact ionization (EI). However, the ESI process is imperfectly “soft” and does produce some degree of in-source fragmentation. Furthermore, secondary adducts, multimers, and fragmentation products of these can form during the ionization process, resulting in multiple observed ions representative of a single compound. These redundant signals are effectively utilized for EI spectra to allow for spectral-matching-based annotation metabolite signals.

Data analysis workflows that seek to detect mass signals in a nontargeted manner utilize both mass and retention time-based specificity—the resulting signal is commonly referred to as a

“feature”. In the absence of co-elution, one feature originates from a single compound. However, the reciprocal is largely untrue: a single compound can give rise to multiple features, as described above. Therefore, many metabolomics data processing tools, including both commercial and open-source tools, attempt to group features into spectra. Some grouping strategies are based on chemically meaningful and predictable patterns reflecting known phenomenon. However, this approach can be compromised by (i) interfering signals from co-eluting metabolites in complex samples that happen to look like fragments, adducts, or isotopes and (ii) unpredictable mass spectral fragments, adducts, or isotopes. As such, an unsupervised approach to grouping features is an attractive alternative. Previous tools including CAMERA,<sup>2</sup> AMDIS,<sup>3</sup> and MSClust<sup>4</sup> have attempted to address this issue, but none of these make full use of the nontargeted data. For example, CAMERA is biased toward the most abundant features and utilizes discrete binning by retention time. MSClust also looks for co-eluting and co-varying features and ultimately selects a representative “centrotype” feature for downstream statistical analysis—the majority of features are discarded. AMDIS works on a single data file, is generally not used for quantitation, and

**Received:** April 25, 2014

**Accepted:** June 13, 2014

73 does not utilize high-mass-accuracy data. Furthermore, all of  
74 these tools are designed for single-channel MS datasets.

75 Here, we report the development of a novel metabolomics  
76 workflow constructed around indiscriminant MS/MS (idMS/  
77 MS) data acquisition, which employs high-collision-energy  
78 fragmentation without precursor ion selection,<sup>5</sup> acquired  
79 concurrently with low-collision-energy MS data. Our method  
80 is based on the premise that two features resulting from the  
81 same compound exhibit similarity in their retention times and a  
82 high correlation in their abundance profiles across different  
83 samples within a dataset. Based on this observation, we have  
84 developed a simple similarity function between features that  
85 allows us to use hierarchical clustering to generate the spectra  
86 of chemical compounds by grouping features from a single  
87 compound in a single cluster. Feature finding is conducted in  
88 both low- and high-collision-energy data, and a custom feature  
89 similarity score drives clustering of features into spectra suitable  
90 for informed manual interpretation, as well as automated  
91 database searching. This approach results in both in-source MS  
92 and idMS/MS spectra for all detected features and enables  
93 spectral matching to public, commercial, and custom spectral  
94 databases without additional experimentation.

## 95 ■ EXPERIMENTAL SECTION

96 **Sample Acquisition and Preparation.** Equine cerebro-  
97 spinal fluid (CSF) samples were obtained as previously  
98 described.<sup>6</sup> CSF was thawed at 4 °C, and 100 μL of CSF was  
99 precipitated with 400 μL of cold methanol. This solution was  
100 mixed thoroughly, incubated at -20 °C for 1 h, and spun at 12  
101 000g for 15 min to remove proteins. The supernatant was  
102 transferred to autosampler vials for UPLC-MS analysis. The  
103 validation dataset consists of 50 urine samples, collected from  
104 Swedish males. Samples were prepared by thawing the urine at  
105 4 °C, diluting with equal parts water, and centrifuging to  
106 remove particulates.

107 **UPLC-MS Data Acquisition.** Metabolome analysis of CSF  
108 and urine samples were accomplished using a Waters Acquity  
109 UPLC system coupled to a time-of-flight mass spectrometer  
110 (Xevo G2 Q-TOF MS). Five microliters (5 μL) of either  
111 protein-depleted CSF or diluted urine was injected onto an  
112 HSS T3 column (Waters, 1 mm × 100 mm, 1.7 μM), and  
113 eluted using a gradient of water to acetonitrile, each containing  
114 0.1% formic acid. The gradient was held at 0.1% B for 1 min,  
115 ramped to 95% B over 12 min, and held for 3 min, before  
116 returning to 0.1% B and equilibrating for 3.9 min (20 min run  
117 time). The flow rate was held constant at 200 μL/min. Eluent  
118 was ionized via positive-mode electrospray ionization, with  
119 capillary voltage set to 2.2 kV, cone to 30 V, extraction cone to  
120 2, with a source temperature of 150 °C and the desolvation  
121 nitrogen gas set to 350 °C at a flow rate of 800 L/h. Before  
122 acquisition, the instrument was calibrated via an infusion of  
123 sodium formate to within an error of 1 ppm. Mass accuracy was  
124 ensured via infusion of leucine enkephalin lockmass, collected  
125 as a 0.5 s scan at a collision energy of 10 V every 20 s. Sample  
126 data were acquired in MSE mode, with alternating scans (0.2 s/  
127 scan, *m/z* 50–1200) collected at collision energy of 6 V (MS)  
128 or using a CE ramp from 15 V to 30 V (idMS/MS). Each

sample was injected in duplicate, with each set of injections  
being completely randomized for acquisition order. In addition,  
the samples were analyzed using data-dependent acquisition  
mode for traditional MS/MS experiments, with one DDA MS/  
MS spectrum acquired per MS scan, with a minimum precursor  
intensity threshold of 200 counts per second. All data were  
acquired in centroid mode.

**Raw Data Conversion and Processing.** Waters raw files  
were converted to cdf format using Databridge, which separates  
low-collision-energy MS and high-collision-energy idMS/MS  
data into two separate cdf files. The lockmass function data was  
discarded for this application. Feature detection (utilizing the  
centWave algorithm), an initial grouping step using a wide  
bandwidth (3), retention time correction, regrouping using a  
narrow bandwidth (1.5), and peak filling was performed using  
XCMS<sup>7</sup> (v. 1.32.0) in R<sup>8</sup> (v. 2.15). CAMERA<sup>2</sup> (v. 1.16.0) was  
used a benchmark comparison, utilizing default values.

**RAMClust Approach.** The RAMClust approach was  
developed in Matlab and is currently fully implemented in R  
in a package called RAMClustR, and it is currently available via  
github (<https://github.com/cbroeckl/RAMClustR>). Imple-  
mentation in R allowed an XCMS object to be used directly  
as input. The data within the XCMS object were extracted  
using the XCMS groupval function and was normalized to the  
total XCMS extracted ion signal (the quantile<sup>9</sup> method is an  
available option in RAMClustR). When a second collision  
energy level is used (as is possible with Waters M<sup>SE</sup> datasets  
utilized in this study), the user directs delineation of MS and  
idMSMS datasets using a tag located within the filename or  
filepath of the xcms object. RAMClustR is also capable of  
accepting properly formatted data matrices from other peak  
detection tools, with the only requirements being:

- (1) no more than one sample (or file) name column and one  
feature name row;
- (2) feature names that contain the mass and retention times,  
separated by a constant delimiter; and
- (3) features in columns and samples in rows.

If both MS and idMS/MS data are to be imported, the feature  
names must be identical between the two datasets.

RAMClust similarity was calculated for the full feature matrix  
(within a user-specified maximum-allowed retention time  
window). Metabolomics datasets can generate thousands to  
tens of thousands of features, which can tax the memory of  
many desktop computers. To manage memory, we utilize the ff  
package,<sup>10</sup> which allows for rapid temporary storage of large R  
objects using physical disk space rather than in memory, and  
process large data matrices in square blocks (2000 features at a  
time by default). The RAMClust similarity scoring utilizes a  
Gaussian function, allowing flexibility in tuning correlational  
and retention time similarity decay rates independently, based  
on the dataset and the acquisition instrumentation. The  
correlational relationship between two features can be  
described by either MS-MS, MS-idMS/MS, or idMS/MS-  
idMS/MS values, and we use Pearson's correlation to calculate  
similarity:

$$S_{ij} = \max \left\{ \exp \left[ -\frac{(1 - c_{ij}^{MS1/MS1})^2}{2\sigma_1^2} \right], \exp \left[ -\frac{(1 - c_{ij}^{MS2/MS2})^2}{2\sigma_2^2} \right], \exp \left[ -\frac{(1 - c_{ij}^{MS1/MS2})^2}{2\sigma_{12}^2} \right] \right\} \exp \left[ -\frac{(t_i - t_j)^2}{2\sigma_t^2} \right]$$

185 where  $(c_{ij}^{MS1/MS2})'$  is the correlation coefficient between  $x_i^{MS1}$   
186 and  $x_j^{MS2}$  ( $i$  and  $j$  represent the peak areas in each sample for  
187 any two features), and  $\sigma_t$  and  $\sigma_r$  represent sigma values for the  
188 retention time and correlational  $r$  value, respectively.  
189 Similarities were then converted to dissimilarities ( $D_{ij} = 1 -$   
190  $S_{ij}$ ) for clustering. The output similarity matrix was then  
191 clustered using average (for this study) or complete linkage  
192 hierarchical clustering via that fastcluster package.<sup>11</sup> The  
193 dendrogram was then cut using the cutreeDynamicTree  
194 function in the package, dynamicTreeCut.<sup>12</sup> For this  
195 application, the minimum module size is set to 2, dictating  
196 that only clusters with two or more features are returned, as  
197 singletons are impossible to interpret intelligently.

198 Cluster membership, in conjunction with the abundance  
199 values from individual features in the input data, were used to  
200 create spectra. Mass was derived from the feature mass, and the  
201 abundance for each mass in the spectrum was derived from the  
202 weighted mean of the intensity values for that feature. These  
203 spectra were then exported as an msp formatted document,  
204 which can be directly imported by NIST MSsearch, or used as  
205 input for MassBank<sup>13</sup> or NIST msPepSearch ([http://peptide.nist.gov/software/ms\\_pep\\_search\\_gui/MSPepSearch.html](http://peptide.nist.gov/software/ms_pep_search_gui/MSPepSearch.html))  
206 batch searching. Finally, the cluster membership was then used  
207 to create a third dataset, SpecData, which represented the MS  
208 level data after condensing feature intensities into spectral  
209 intensities using a weighted mean function, where the more-  
210 abundant signals contribute more to the spectral intensity.

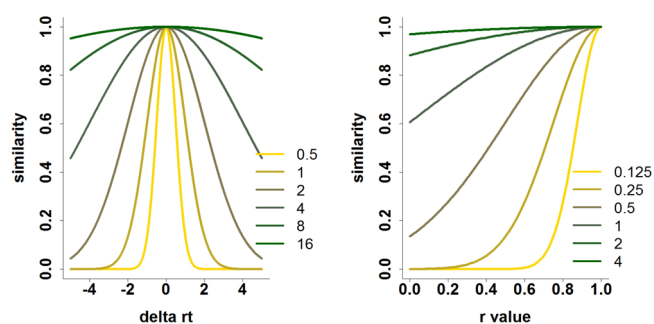
## 212 ■ RESULTS AND DISCUSSION

213 We developed and tested our approach using a UPLC-MS  
214 dataset of 38 samples of equine cerebrospinal fluid, and  
215 subsequently validated the approach in an independent urine  
216 dataset (see Figure 1 in the supplementary material). XCMS<sup>7</sup>  
217 was used for feature finding, retention time correction, and  
218 alignment, and the resulting dataset was subsequently  
219 normalized to total XCMS signal intensity for each sample.  
220 The output data was then divided into low-collision-energy  
221 (MS) and high-collision-energy (idMS/MS) datasets, each with  
222 dimensions of row number equal to the number of injections  
223 and column number equal to the number of features (21060,  
224 for the CSF dataset). Each cell of these datasets represents the  
225 signal intensity at either low (MS) or high (idMS/MS) collision  
226 energy. We developed a custom similarity matrix, which is the  
227 product of two Gaussian terms: one that considers the  
228 differences in retention times between two features and a  
229 second that considers the correlation between two features  
230 across all samples in the dataset. These two terms have widths  
231 defined by  $\sigma_t$  and  $\sigma_r$ , respectively. This captures our intuition  
232 that two features are similar if they are close in retention time  
233 and are correlated: both are required for two features to be  
234 grouped. Following the computation of the similarity matrix,  
235 features are clustered using hierarchical clustering.

236 To generate discrete clusters from the resulting hierarchical  
237 clustering dendrogram, we then used the DynamicTreeCut<sup>14</sup>  
238 package in R. Cluster membership of each feature provides  
239 qualitative spectral membership information, and the quantitative  
240 data are taken from the MS and idMS/MS datasets;  
241 abundance values are calculated as the averaged signal intensity  
242 for each feature separately in both the low- and high-collision-  
243 energy datasets. Thus, for each cluster, two spectra are re-  
244 created, corresponding to the low-collision-energy in-source  
245 spectra and the high-collision-energy counterparts.

Any feature clustering tool must demonstrate accuracy to be  
246 useful in reducing redundancy without reducing biological  
247 coverage. One option to accomplish this is to compare the  
248 results of the clustering to a small panel of known compounds  
249 that are spiked into a sample. While this is a valid approach, it  
250 relies on the assumption that the chosen panel of compounds is  
251 representative of all the metabolites in a complex biological  
252 matrix. Thus, to increase the breadth of our validation  
253 experiments, we instead assessed the accuracy of the clustering  
254 by comparison against MS/MS spectra acquired using a  
255 traditional dependent acquisition (DDA) approach from the  
256 same CSF samples. All precursor ions that (i) could be mapped  
257 to a feature in the output dataset and (ii) contained more than  
258 10 product ions were used as “valid” spectra for comparison.  
259 These spectra represented known precursor-product ion  
260 relationships from many of the major signals in the dataset,  
261 even if the identity of the compounds was unknown. The  
262 spectra created by RAMClust were then compared to the DDA  
263 spectra and the dot product spectral similarity score was  
264 calculated as a measure of accuracy, as described previously.<sup>15</sup>  
265 While the complexity of in-source and indiscriminant MS/MS  
266 signals is expected to be higher than DDA MS/MS spectra for  
267 the same compound, more-accurate clustering will still be  
268 revealed as relatively higher dot-product similarity scores  
269 between the RAMclustR reconstructed spectra and the mapped  
270 DDA MS/MS spectrum.

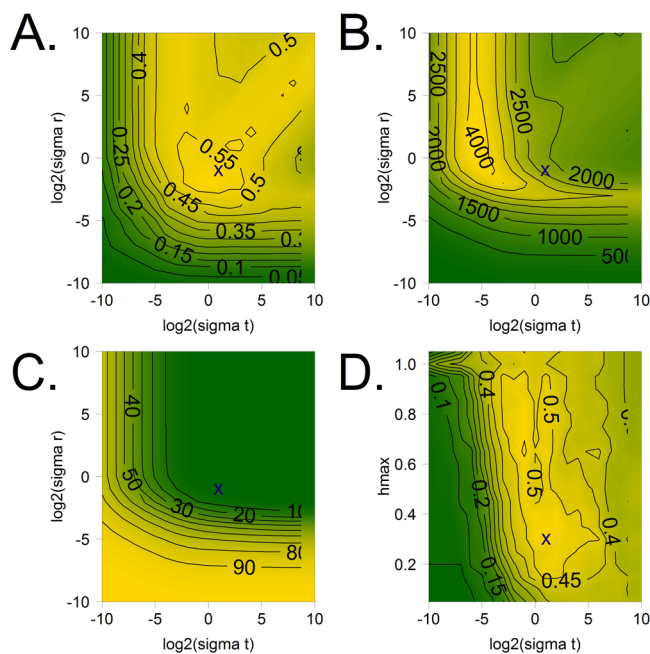
271 The RAMClust algorithm has several parameters that can be  
272 tuned by the user to improve clustering accuracy. Parameters  $\sigma_t$   
273 and  $\sigma_r$  represent Gaussian tuning parameters of retention time  
274 similarity and correlational score, respectively, between feature  
275 pairs. The influence of these two parameters on the similarity is  
276 depicted in Figure 1. These tuning parameters will allow the  
277



**Figure 1.** RAMClust is based on a custom feature similarity score, which is the product of two terms that capture similarity in retention time and correlation across samples. Each of the two terms has a tuning parameter associated with it that controls the width of the corresponding Gaussian:  $\sigma_t$  for retention time (left) and  $\sigma_r$  for the degree of (right). Increased values for the two  $\sigma$  terms decrease the rate of decay in the similarity score, as a function of either retention time difference or correlation  $r$  between pair of features.

algorithm to be used with MS data from any chromatographic  
278 platform. When idMS/MS data are available, correlational  
279 similarity can be calculated between two features, at the level of  
280 either MS vs MS, MS vs idMS/MS, or idMS/MS vs idMS/MS.  
281 While the MS-idMSMS correlation theoretically represents the  
282 CID event most directly, this relationship is subject to potential  
283 interfering signals in both data channels (MS and idMS/MS).  
284 In practice, a strong correlational relationship at any of the  
285 three levels represents strong evidence of precursor-product  
286

287 relationships; thus, the algorithm utilized the maximum  
 288 correlational  $r$ -value of the three relationships.  
 289 The influence of  $\sigma_t$  and  $\sigma_r$  on the average spectral similarity  
 290 between RAMClust and DDA spectra was rigorously evaluated  
 291 at 441 combinations of parameter levels of  $\sigma_t$  and  $\sigma_r$  (Figure  
 292 2a). These results revealed a plateau of high spectral similarity



**Figure 2.** Influence on RAMClustR parameters,  $\sigma_t$  for  $\sigma_t$  for and  $h_{max}$  were systematically varied to examine the influence of these parameters on feature grouping accuracy, the number of clusters, and the number of ungrouped features (singletons). (A) RAMClust spectra generated using  $\sigma_t$  and  $\sigma_r$  values of 2 and 0.5 produce the strongest dot product similarity to DDA spectra, which represent validated precursor product relationships. This  $\sigma_t$  value is roughly half the median XCMS peak width, indicating that the  $\sigma_t$  value can be set automatically when XCMS data are used as the input. (B) Influence of  $\sigma_t$  and  $\sigma_r$  on the number of clusters with at least two features. The optimal values  $\sigma_t$  and  $\sigma_r$  (denoted with an “x”), as determined by the maximal dot product similarity, results in  $\sim 2500$  clusters. (C)  $\sigma_t$  and  $\sigma_r$  values that are too selective results in fewer clusters, because of high singleton (features which cluster with no other features). (D) The dot product similarity scoring benefits from some precutting of the tree, as provided by the dynamicTreeCut algorithm, allowing us to set a default maximal cluster height of 0.3.

293 at values of  $\sigma_t = 2$  and  $\sigma_r = 0.5$  (Figure 2a). This  $\sigma_t$  value was  
 294 approximately half the median peak width of the XCMS  
 295 detected peak (max-min time for each individual peak in the  
 296 xcms object), indicating that we can directly use XCMS input  
 297 to set this parameter without user intervention: this holds true  
 298 for an independent dataset of urine samples (see the  
 299 supplementary material). Correlation is a scale-free statistic,  
 300 and it should be platform-neutral; thus, we used our observed  
 301 optimal value of 0.5 and can expect reasonable performance on  
 302 any platform. Implementation of RAMClustR using parameters  
 303 that maximized MS/MS similarity between reconstructed  
 304 spectra and DDA spectra generated  $\sim 2500$  clusters with at  
 305 least two features (Figure 2b), and relatively few singletons  
 306 (Figure 2c). This algorithm generated a large stable region,  
 307 indicating that it is robust to small changes in parameter values.  
 308 This stability generated a strong MS/MS similarity, even at

“unreasonable”  $\sigma_t$  values ( $>200$  s), as long as  $\sigma_r$  is proportion- 309  
 ally high (Figure 2a). We interpret this as a scaling 310  
 phenomenon, as the dynamicTreeCut algorithm is responsive 311  
 to tree “shape” rather than an absolute height.<sup>14</sup> The 312  
 dynamicTreeCut maximum height parameter was also exam- 313  
 ined in conjunction with  $\sigma_r$ , and it revealed that the tree 314  
 pruning step benefited from some precutting (Figure 2d); thus, 315  
 we employ a default value of 0.3 for this parameter. These 316  
 parametrization rules make the algorithm extremely easy to use: 317  
 when an XCMS object is used as input, the user needs to set 318  
 none of these parameters, and when a dataset is imported from 319  
 other software, only  $\sigma_t$  needs to be manually set. The output 320  
 MS/MS similarity using default RAMClust similarity scores was 321  
 used to compare results against the only other feature grouping 322  
 tool in R: CAMERA. The results of this comparison indicated 323  
 that RAMClust grouping of features resulted in spectra that are 324  
 more similar to DDA spectra than the results generated from 325  
 CAMERA’s groupFWHM, groupCorr, and groupDen functions 326  
 (see Table 1). This observation was validated on a second LC- 327 11

**Table 1. Comparison between RAMClustR and CAMERA<sup>a</sup>**

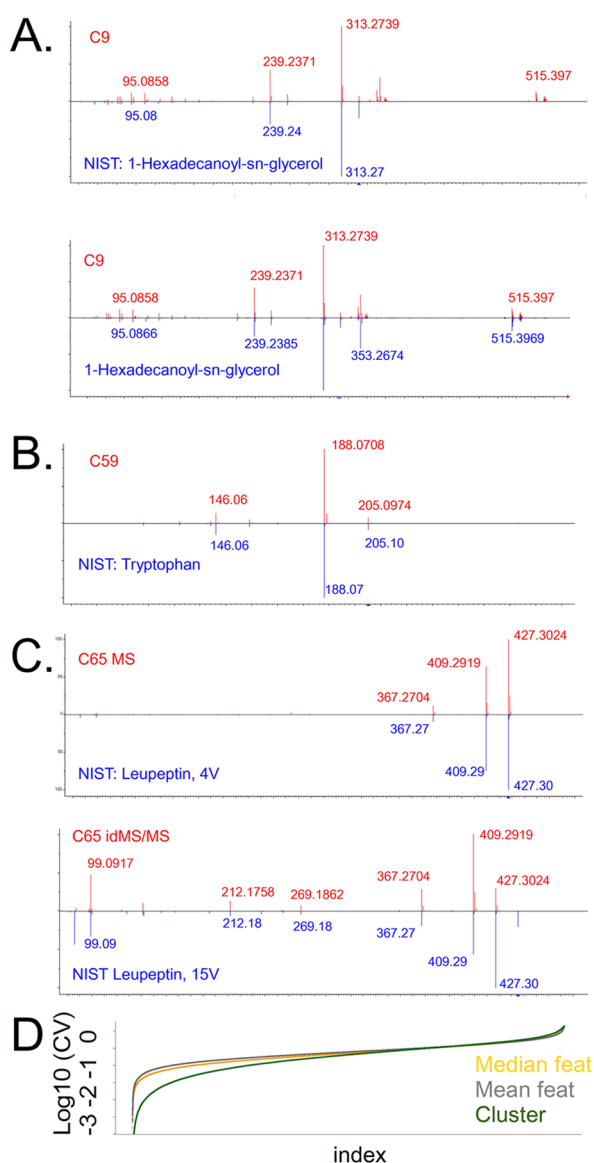
method	MSMS similarity <sup>b</sup>	nClus (>1) <sup>c</sup>	perSing <sup>d</sup>
<b>CSF Dataset</b>			
xsb $\leftarrow$ groupFWHM(xset)	0.202	535	0.43
xsc $\leftarrow$ groupCorr(xsb)	0.177	784	29.56
xsd $\leftarrow$ groupDen(xsa)	0.043	39	0.00
RAMClustR(xset)	0.382	3248	15.47
<b>Urine Dataset</b>			
xsb $\leftarrow$ groupFWHM(xset)	0.106	290	4.88
xsc $\leftarrow$ groupCorr(xsb)	0.059	332	61.77
xsd $\leftarrow$ groupDen(xsa)	0.020	35	0.00
RAMClustR(xset)	0.228	827	32.75

<sup>a</sup>The comparisons were performed using default values for both the CSF and Urine datasets. The first three rows in both the CSF and Urine datasets reflect CAMERA functions, while the final row reflects RAMClustR-based grouping. <sup>b</sup>MSMSsimilarity refers to the spectral similarity between mapped feature for which data-dependent MS/MS data were available and the reconstructed spectra from the output dataset defined in the “method” column. <sup>c</sup>nClus (>1) refers to the number of clusters with two or more features defined by the grouping method. <sup>d</sup>perSing is the percentage of all features in the data set that remain ungrouped (singletons).

MS dataset of urine samples: RAMClust grouping resulted in 328  
 clustering output that better represents valid feature relation- 329  
 ships and, consequentially, biological small molecule signals. 330

The spectra produced via RAMClust grouping can written to 331  
 NIST MSP format for viewing and searching, and they can be 332  
 submitted directly to the MassBank Database<sup>13</sup> batch search 333  
 tool, submitted for batch searching to NIST msPepSearch, 334  
 and/or viewed and searched via the NIST MSSearch program. 335  
 All these tools offer the ability to generate and search against 336  
 custom libraries of spectra, and our laboratory is creating 337  
 libraries of in-source spectra toward this end. However, idMS/ 338  
 MS spectra re-created from the RAMClust algorithm and 339  
 workflow were highly similar to authentic NIST MS/MS 340  
 database spectra (see Figures 3a–c), demonstrating that this 341 13  
 workflow can take full advantage of existing resources. 342

Since RAMClust-generated spectra accurately reflect spectra 343  
 of authentic chemical standards, the intensity of the spectra 344  
 themselves can be used as the quantitative unit for downstream 345  
 statistical analysis. The intensity of the spectra were calculated 346  
 using a weighted mean function of all the component features, 347



**Figure 3.** (A) Cluster membership and peak area data are used to generate spectra, which can be searched against spectra databases. The in-source low-collision-energy spectrum representing C9 was identified as hexadecanoyl-*sn*-glycerol (16:0 MAG) in the CSF samples, and shows a strong match to the NIST library spectrum representing this compound. However, the match is even stronger if all the in-source phenomenon are considered (bottom panel, standard run by the authors under identical analytical conditions). (B) Tryptophan in-source low-collision-energy spectrum can be identified with a high degree of confidence from either NIST MS/MS spectra (top) or a custom library spectrum (bottom). (C) Both low-collision-energy spectra (top) and high-collision-energy spectra (bottom) can be used for the same compound to increase the confidence of identification in the event that the MS spectrum is sparse, as demonstrated by leupeptin, a protease inhibitor added to the CSF samples before processing. (D) Clustering of features results in reduced analytical variation. The coefficient of variation (CV) of all individual compound measurements was calculated for all clusters, and compared to the median or mean feature CV for the features comprising those clusters. These ~120 000 measures of variation indicate that the analytical variation for the majority of compound measurements is greatly reduced through aggregation into compound clusters or spectra.

sample in the dataset. The use of spectra dramatically reduced analytical variation through an averaging of measurement noise, as compared to either the mean or median feature-based variation for each cluster (see Figure 3D).

## CONCLUSIONS

Annotation of mass signals in nontargeted metabolomics experiments remains a significant bottleneck and is arguably one of the most important challenges to the field as confident metabolite identification is required for biological interpretation. In this report, we demonstrate a novel workflow utilizing indiscriminant MS/MS data acquisition, expanded feature finding and a novel clustering algorithm to group features based on both low- and high-collision-energy data to generate spectra that are compatible with publically available spectral search tools. The workflow allows for more-efficient use of instrumentation, reduced feature redundancy and false discovery rate correction burden for downstream univariate statistical tests, improved analytical reproducibility, a more-automated annotation workflow, and greatly increased confidence in the annotations, compared to accurate mass-based searching alone. RAMClustR is available for download at <https://github.com/cbroeckl/RAMClustR>.

## AUTHOR INFORMATION

### Corresponding Authors

\*Tel.: 970-491-2273. E-mail: Corey.Broeckling@colostate.edu (C. D. Broeckling).

\*Tel.: 970-213-9093. E-mail: afsar@rams.colostate.edu (Fayyaz ul Amir Afsar Minhas).

\*E-mail: sneumann@ipb-halle.de (Steffen Neumann).

\*Tel.: 970-491-4068. E-mail: asa@cs.colostate.edu (Asa Ben-Hur).

\*Tel.: 970-491-0961. E-mail: Jessica.Prenni@colostate.edu (Jessica E. Prenni).

### Notes

The authors declare no competing financial interest.

## ACKNOWLEDGMENTS

F.A.A. was funded by the Fulbright Scholarship Program of the U.S. Department of State and the Higher Education Commission of the Government of Pakistan.

## REFERENCES

- Whitehouse, C. M.; Dreyer, R. N.; Yamashita, M.; Fenn, J. B. Electrospray Interface for Liquid Chromatographs and Mass Spectrometers. *Anal. Chem.* **1985**, *57*, 675–679 (DOI: 10.1021/ac00280a023).
- Kuhl, C.; Tautenhahn, R.; Böttcher, C.; Larson, T. R.; Neumann, S. CAMERA: An Integrated Strategy for Compound Spectra Extraction and Annotation of Liquid Chromatography/Mass Spectrometry Data Sets. *Anal. Chem.* **2012**, *84*, 283–289 (DOI: 10.1021/ac202450g).
- Halket, J. M.; Przyborowska, A.; Stein, S. E.; Mallard, W. G.; Down, S.; Chalmers, R. A. Deconvolution gas chromatography mass spectrometry of urinary organic acids—Potential for pattern recognition and automated identification of metabolic disorders. *Rapid Commun. Mass Spectrom.* **1999**, *13*, 279–284 (DOI: 10.1002/(SICI)1097-0231(19990228)13:4<279::AID-RCM478>3.0.CO;2-1).
- Tikunov, Y. M.; Laptinok, S.; Hall, R. D.; Bovy, A.; de Vos, R. C. M. MSclust: A tool for unsupervised mass spectra extraction of chromatography–mass spectrometry ion-wise aligned data. *Metabolomics* **2012**, *8*, 714–718 (DOI: 10.1007/s11306-011-0368-2).

such that each value in the resulting dataset represents the quantitative signal intensity value for each spectrum for each

- 409 (5) (a) Plumb, R. S.; Johnson, K. A.; Rainville, P.; Smith, B. W.;  
410 Wilson, I. D.; Castro-Perez, J. M.; Nicholson, J. K. UPLC/MS<sup>E</sup>: A new  
411 approach for generating molecular fragment information for biomarker  
412 structure elucidation. *Rapid Commun. Mass Spectrom.* **2006**, *20*, 1989–  
413 1994 (DOI: 10.1002/rcm.2550). (b) Plumb, R. S.; Johnson, K. A.;  
414 Rainville, P.; Smith, B. W.; Wilson, I. D.; Castro-Perez, J. M.;  
415 Nicholson, J. K. UPLC/MS<sup>E</sup>: A new approach for generating  
416 molecular fragment information for biomarker structure elucidation  
417 (Erratum). *Rapid Commun. Mass Spectrom.* **2006**, *20*, 2234 (DOI:  
418 10.1002/rcm.2602).
- 419 (6) Broccardo, C. J.; Hussey, G. S.; Goehring, L.; Lunn, P.; Prenni, J.  
420 E. Proteomic Characterization of Equine Cerebrospinal Fluid. *J. Equine*  
421 *Vet. Sci.* **2013**, <http://dx.doi.org/10.1016/j.jevs.2013.07.013>.
- 422 (7) Smith, C. A.; Want, E. J.; O'Maille, G.; Abagyan, R.; Siuzdak, G.  
423 XCMS: Processing Mass Spectrometry Data for Metabolite Profiling  
424 Using Nonlinear Peak Alignment, Matching, And Identification. *Anal.*  
425 *Chem.* **2006**, *78*, 779–787 (DOI: 10.1021/ac051437y).
- 426 (8) Team, R. C. Vienna, Austria, 2013.
- 427 (9) (a) Bolstad, B. M. *preprocessCore: A collection of pre-processing*  
428 *functions, R package*, Version 1.22.0, 2014. (b) Brodsky, L.; Moussaieff,  
429 A.; Shahaf, N.; Aharoni, A.; Rogachev, I. Evaluation of Peak Picking  
430 Quality in LC-MS Metabolomics Data. *Anal. Chem.* **2010**, *82*, 9177–  
431 9187 (DOI: 10.1021/ac101216e).
- 432 (10) Adler, D.; Gläser, C.; Nenadic, O.; Oehlschlägel, J.; Zucchini, W.  
433 *R package version 2.2*, 11th Edition, 2013.
- 434 (11) Mullner, D. *J. Stat. Software* **2013**, *53*, 1–18.
- 435 (12) Langfelder, P.; Zhang, B.; Horvath, S. Defining clusters from a  
436 hierarchical cluster tree: The Dynamic Tree Cut package for R.  
437 *Bioinformatics* **2008**, *24*, 719–720 (DOI: 10.1093/bioinformatics/  
438 btm563).
- 439 (13) Horai, H.; Arita, M.; Kanaya, S.; Nihei, Y.; Ikeda, T.; Suwa, K.;  
440 Ojima, Y.; Tanaka, K.; Tanaka, S.; Aoshima, K.; Oda, Y.; Kakazu, Y.;  
441 Kusano, M.; Tohge, T.; Matsuda, F.; Sawada, Y.; Hirai, M. Y.;  
442 Nakanishi, H.; Ikeda, K.; Akimoto, N.; Maoka, T.; Takahashi, H.; Ara,  
443 T.; Sakurai, N.; Suzuki, H.; Shibata, D.; Neumann, S.; Iida, T.; Tanaka,  
444 K.; Funatsu, K.; Matsuura, F.; Soga, T.; Taguchi, R.; Saito, K.;  
445 Nishioka, T. MassBank: A public repository for sharing mass spectral  
446 data for life sciences. *J. Mass Spectrom.* **2010**, *45*, 703–714 (DOI:  
447 10.1002/Jms.1777).
- 448 (14) Langfelder, P.; Zhang, B.; Horvath, S. Defining clusters from a  
449 hierarchical cluster tree: The Dynamic Tree Cut package for R.  
450 *Bioinformatics* **2008**, *24*, 719–720 (DOI: DOI 10.1093/bioinfor-  
451 matics/btm563).
- 452 (15) Broeckling, C. D.; Heuberger, A. L.; Prince, J. A.; Ingelsson, E.;  
453 Prenni, J. E. Assigning precursor–product ion relationships in  
454 indiscriminant MS/MS data from non-targeted metabolite profiling  
455 studies. *Metabolomics* **2013**, *9*, 33–43 (DOI: 10.1007/s11306-012-  
456 0426-4).