1  # Chemically-informed Analyses of Metabolomics Mass Spectrometry Data

2  # with Qemistree

3  Authors: Anupriya Tripathi[1,2,3#], Yoshiki Vázquez-Baeza[4,5#], Julia M. Gauglitz[3,6], Mingxun Wang[3], Kai Dührkop[7],

4  Mélissa Nothias-Esposito[3], Deepa D. Acharya[3,8], Madeleine Ernst[3,6,9], Justin J.J. van der Hooft[10], Qiyun Zhu[2],

5  Daniel McDonald[2], Antonio Gonzalez[2], Jo Handelsman[8], Markus Fleischauer[7], Marcus Ludwig[7], Sebastian Böcker[7],

6  Louis-Félix Nothias[3], Rob Knight[2,4,5,11], Pieter C. Dorrestein[3,5,6]

7

8  Author Affiliations:

9  [1] Division of Biological Sciences, University of California San Diego.

10  [2] Department of Pediatrics, University of California San Diego.

11  [3] Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California San Diego

12  [4] Jacobs School of Engineering, University of California San Diego, La Jolla, California, USA.

13  [5] Center for Microbiome Innovation, University of California San Diego, La Jolla, California, USA.

14  [6] Collaborative Mass Spectrometry Innovation Center, Skaggs School of Pharmacy and Pharmaceutical Sciences,

15  University of California, San Diego, La Jolla, CA, USA.

16  [7] Chair for Bioinformatics, Friedrich-Schiller-University, Jena, Germany.

17  [8] Wisconsin Institute of Discovery, University of Wisconsin-Madison, Madison, Wisconsin, USA

18  [9] Section for Clinical Mass Spectrometry, Department of Congenital Disorders, Danish Center for Neonatal

19  Screening, Statens Serum Institut, Copenhagen, Denmark

20  [10] Bioinformatics Group, Plant Sciences Group, Wageningen University, Wageningen, The Netherlands.

21  [11] Department of Computer Science and Engineering, University of California San Diego.

22  [#] Equal contribution.

23

24  Author contributions:

25  PCD, AT conceived the concept and managed the project.

26  AT and YVB developed the algorithm and wrote the code for Qemistree.

27  AT and YVB contributed equally to the work.

28  LFN, RK, PCD supervised method implementation.

29  KD, MW, JJJvdH, ME, DM, and AG tested and provided suggestions on how to improve the method.

30  MW managed the deployment of Qemistree on GNPS.

31  AT and MW developed the GNPS-Qemistree Dashboard.

32  DA and AT wrote the documentation for the GNPS-Qemistree workflow.

33  YVB, QZ, and AT developed Qemistree-iTOL visualization.

34  LFN and MNE performed the mass-spectrometry for the evaluation dataset.

35  AT, YVB, and LFN analyzed and interpreted the evaluation data.

36  JMG performed mass spectrometry of the Global Foodomics samples.

37  AT, JMG analyzed and interpreted the Global Foodomics data.

38    KD, MF, ML, and SB supported the integration of SIRIUS, Zodiac, and CSI:FingerID.

39    PCD, AT, YVB, and RK wrote the manuscript.

40    LFN, JMG, MNE, JJJvdH, ME, KD, QZ, DM, AG, JH, MF, ML, SB, and RK improved the manuscript.

## Abstract

42    Untargeted mass spectrometry is employed to detect small molecules in complex biospecimens,

43    generating data that are difficult to interpret. We developed Qemistree, a data exploration

44    strategy based on hierarchical organization of molecular fingerprints predicted from

45    fragmentation spectra, represented in the context of sample metadata and chemical ontologies.

46    By expressing molecular relationships as a tree, we can apply ecological tools, designed around

47    the relatedness of DNA sequences, to study chemical composition.

## Main

49    Molecular networking[1], introduced in 2012, was one of the first data organization approaches to

50    visualize the relationships between fragmentation spectra for similar molecules from tandem

51    mass spectrometry data in the context of metadata. It formed the basis for the web-based mass

52    spectrometry infrastructure, Global Natural Products Social Molecular Networking[2] (GNPS,

53    https://gnps.ucsd.edu/) which sees ~200,000 new accessions per month. Molecular networking is

54    used for a range of applications[3] in drug discovery, environmental monitoring, medicine, and

55    agriculture. While molecular networking is useful for visualizing closely related molecular

56    families, the inference of chemical relationships at a dataset-wide level and in the context of

57    diverse metadata requires complementary representation strategies. To address this need, we

58    developed an approach that uses fragmentation trees[4] and supervised machine learning[5] to

59    calculate all pairwise chemical relationships and visualizes it in the context of sample metadata

60    and molecular annotations. We show that a chemical tree enables the application of various tree-

61    based tools, originally developed for analyzing DNA sequencing data[6–9], for exploring mass-

62    spectrometry data.

63

64    We introduce Qemistree, pronounced *chemis-tree*, a software that constructs a chemical tree

65    from fragmentation spectra based on predicted molecular fingerprints[10]. Molecular fingerprints

66    are vectors where each position encodes a substructural property of the molecule. Recent

67    methods allow us to predict molecular fingerprints from tandem mass spectra[11–15]. In Qemistree,

68  we use SIRIUS[16] and CSI:FingerID[13] to obtain predicted molecular fingerprints. The users first

69  perform feature detection[17,18] to generate a list of observed ions, referred to as chemical features

70  henceforth, to be analyzed by Qemistree (Fig. S1). SIRIUS then determines the molecular

71  formula of each feature using the isotope and fragmentation patterns, and estimates the best

72  fragmentation tree explaining the fragmentation spectrum. Subsequently, CSI:FingerID operates

73  on the fragmentation trees using kernel support vector machines to predict molecular properties

74  (2936 properties; Table S1). We use these molecular fingerprints to calculate pairwise distances

75  between chemical features that are hierarchically clustered to generate a tree representing their

76  structural relationships. Although alternative approaches to hierarchically cluster features based

77  on cosine similarity of fragmentation spectra exist[19–21], we use molecular fingerprints as it allows

78  us to compare features based on a diverse range of structural properties predicted by

79  CSI:FingerID. Additionally, as CSI:FingerID was shown to perform well for automatic *in silico*

80  structural annotation[22], we leverage it to search molecular structural databases to provide

81  complementary insights into structures when no match is obtained against spectral libraries.

82  Subsequently, we use ClassyFire[23] to assign a 5-level chemical taxonomy (kingdom, superclass,

83  class, subclass, and direct parent) to all molecules annotated via spectral library matching and *in*

84  *silico* prediction.

85

86  Phylogenetic tools such as iTOL[24] can be used to visualize Qemistree trees interactively in the

87  context of sample information and feature annotations for easy data exploration. The outputs of

88  Qemistree can also be plugged into other workflows in QIIME 2[25] (many of which were

89  originally developed for microbiome sequence analysis) or in R, Python etc. for system-wide

90  metabolomic data analyses [6,7,9, 26]. Qemistree is available to the microbiome community as a

91  QIIME 2 plugin (https://github.com/biocore/q2-qemistree) and the metabolomics community as

92  a workflow on GNPS[2] (https://ccms-ucsd.github.io/GNPSDocumentation/qemistree/). The

93  chemical tree from the GNPS workflow can be explored interactively (e.g.

94  https://qemistree.ucsd.edu/).

95

96  To verify that molecular fingerprint-based trees correctly capture the chemical relationships

97  between molecules, we generated an evaluation dataset with two human fecal samples, a tomato

98  seedling sample, and a human serum sample. Mixtures of these samples were prepared by

99    combining them in gradually increasing proportions to generate a set of diverse but related

100    metabolite profiles and untargeted tandem mass spectrometry was used to profile the chemical

101    composition of these samples. Mass-spectrometry was performed twice using different

102    chromatographic gradients causing a non-uniform retention time shift between the two runs. The

103    data processing of these two experiments leads to the same molecules being detected as different

104    chemical features in downstream analysis. In Figure 1a we highlight how these technical

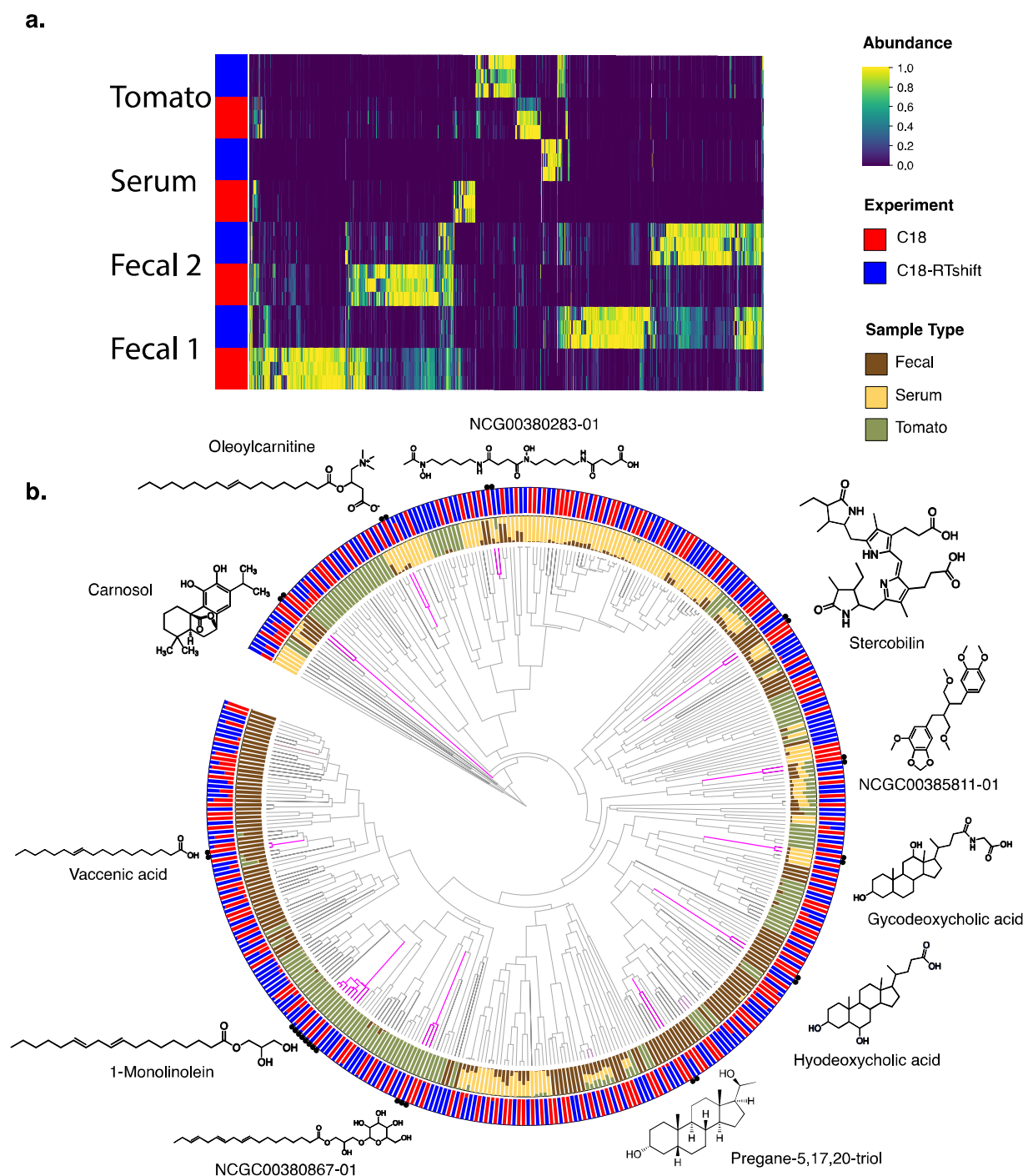105    variations make the same samples appear chemically disjointed.

**Figure 1: Qemistree mitigates aspects of technical artifacts by co-clustering structurally similar molecules across mass spectrometry runs. a)** Sample (y-axis) by molecule (x-axis) heatmap of 2 fecal samples, tomato seedling samples, and serum samples in the evaluation dataset grouped by chromatography conditions. **b)** A chemical tree based on predicted molecular fingerprints representing the structural relationships between compounds detected in the evaluation dataset. Outer ring shows the relative abundance of molecules stratified by mass spectrometry run; inner ring shows the same stratified by fecal, serum and tomato samples in the evaluation dataset. Structurally similar molecules detected as different chemical features due to shift in retention time across mass spectrometry runs are clustered together; we highlight some examples of these artificially duplicated features around the tree. All structures shown are spectral reference library matches obtained from feature-based molecular
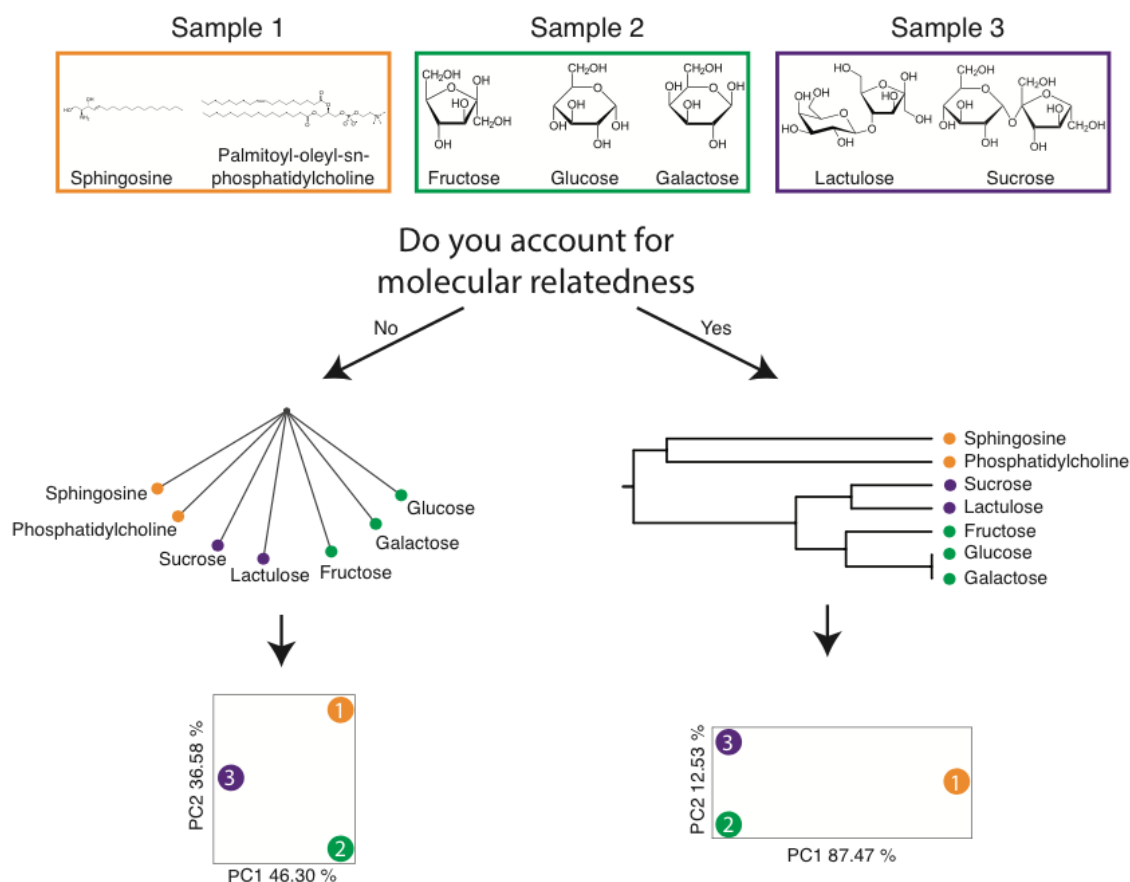
116    networking[17] in GNPS: (https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=efda476c72724b29a91693a108fa5a9d;

117    Metabolomics Standard Initiative (MSI) level 3 annotation)[27].

118

119    Using Qemistree, we map each of the spectra in the two chromatographic conditions (batches) to

120    a molecular fingerprint, and organize these in a tree structure (Fig. 1b). Because molecular

121    fingerprints are independent of retention time shifts, spectra are clustered based on their chemical

122    similarity. This tree structure can be decorated using sample type descriptions, chromatographic

123    conditions, and spectral library matches obtained from molecular networking in GNPS. Figure 1

124    shows that similar chemical features are detected exclusively in one of the two batches.

125    However, based on the molecular fingerprints, these chemical features were arranged as

126    neighboring tips in the tree regardless of the retention time shifts. This result shows how

127    Qemistree can reconcile and facilitate the comparison of datasets acquired on different

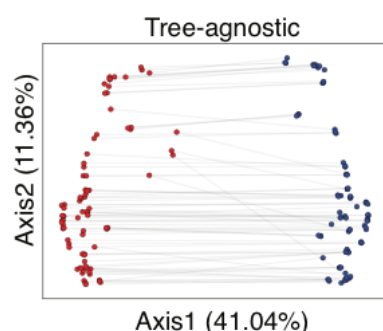128    chromatographic gradients.

129

130    We demonstrate the use of a chemical hierarchy in performing chemically-informed

131    comparisons of metabolomics profiles. In standard metabolomic statistical analyses, each

132    molecule is assumed unrelated to the other molecules in the dataset. Some of the pitfalls of this

133    assumption are highlighted in Figure 2a. Consider a scenario where we want to compare samples

134    1-3. An analysis schema that does not account for the chemical relationships among the

135    molecules in these samples (Figure 2a, left), will assume that the sugars in samples 2 and 3 are as

136    chemically related to the lipids in sample 1 as they are to each other. This would lead to the naive

137    conclusion that samples 1 and 2, and samples 2 and 3 are equally distinct, yet they are not from a

138    chemical perspective. On the other hand, if we account for the fact that sugar molecules are more

139    chemically related to one another than they are to lipids, we can obtain a chemically-informed

140    sample-to-sample comparison. Sedio and coworkers developed the chemical structural

141    compositional similarity (CSCS) metric[28] to account for relationships between molecules based

142    on the similarity of their fragmentation spectra. While CSCS compares samples based on

143    modified cosine scores obtained from molecular networking, we calculate chemical relationships

144    based on structurally-informed molecular fingerprints. We express these relationships in the form

145    of a hierarchy which enables the use of other tree-based tools for downstream data analyses. For

146    example, in Figure 2a, we show that by using a tree of structural relationships between molecules

147    in samples 1-3, we can apply UniFrac[9], a tree-informed distance metric and demonstrate that the

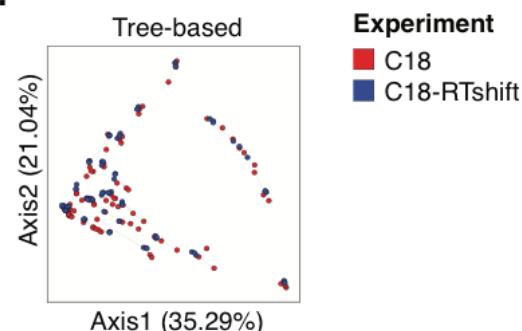148    composition of sample 1 is distinct from samples 2 and 3.

149



150

**Figure 2: The pitfalls of assuming equal relatedness of molecules and the advantages of a chemical tree for**
151
**sample comparison. a)** A scenario where the goal is to compare the chemical composition in samples 1
152
(sphingosine and phosphatidylcholine), 2 (glucose, galactose, and fructose), and 3 (sucrose and lactulose). When we
153
do not account for the chemical relationships between the molecules, i.e. assume that the lipid molecules in sample 1
154
are equally related to the sugars in samples 2 and 3 (left), we conclude that samples 1, 2, and 3 are similarly distinct.
155
If we account for sugar molecules being more chemically related to one another than sugars are to lipid molecules
156

157 (right), we can obtain a chemically-meaningful distance between samples. This is exemplified through a principal
158 coordinates analysis (PCoA) of the computed UniFrac[9] (tree-based) distances among samples; we see that samples 2
159 and 3 are more similar to each other, and sample 1 which is chemically distinct is separated along the primary axis
160 of variation, when distances are computed using the chemical tree. **b, c)** PCoA of samples in the evaluation dataset
161 colored by chromatography conditions. PCoA plot using tree-agnostic (Bray-Curtis[29]) distances which do not
162 account for the chemical relationship between features detected across chromatography conditions (b) and tree-
163 based (Weighted UniFrac[9]) distances which are based on the hierarchical relationships between molecules in the
164 evaluation dataset (c).
165

166 The importance of comparing samples by accounting for their molecular relatedness is

167 highlighted when we contrast the results from ignoring the tree structure (Fig. 2b) to those which

168 integrate it (Fig. 2c). With the structural context provided by Qemistree, the differences between

169 replicates across batches are comparable to the within-batch differences (Fig. S2). The retention

170 time shift in this dataset leads to a strong technical signal that obscures the biological

171 relationships among the samples (permutational ANOVA; tree agnostic[29] pseudo-F=120.75,

172 p=0.001 vs. tree informed[9] pseudo-F=18.2239, p=0.001). We observed and remediated a similar

173 pattern originating from plate-to-plate variation in a recently published study investigating the

174 metabolome and microbiome of captive cheetahs[30] (Fig. S3). In this study, placing the molecules

175 in a tree using Qemistree reduced the observed technical variation (Fig. S3 a, c), and highlighted

176 the dietary effect that was expected (Fig. S3 b, d). These results show how systematic and

177 spurious molecular differences can be mitigated in an unsupervised manner using chemically-

178 informed distance measures based on a tree structure.

179

180 As a case study, we used Qemistree to explore chemical diversity in a set of food samples

181 collected as a part of the Global FoodOmics initiative (http://globalfoodomics.org). We selected

182 a diverse range of food ingredients to represent animal, plant, and fungal groupings[31]. We first

183 performed feature-based molecular networking using MZmine[17,18] to obtain spectral library

184 matches for a subset of the chemical features (~20% annotated with cosine cutoff $> 0.7$).

185 Understanding the chemical relationships between different foods is challenging because most

186 molecules within foods are unannotated. Using Qemistree, we collated GNPS spectral library

187 matches and *in silico* predictions from CSI:FingerID to annotate ~91% of the chemical features

188 (total 634 features after quality filtering) with molecular structures. Using ClassyFire[23], we

189 assigned a chemical taxonomy to 60% of these structures; the remaining 40% returned no

190 ClassyFire taxonomy. Labeling annotations allowed us to retrieve subtrees of distinct chemical

191 classes (Fig. 3a) such as flavonoids, alkaloids, phospholipids, acyl-carnitines, and O-glycosyl

192    compounds in food products. We propagated ClassyFire annotations of chemical features (tree

193    tips) to each internal node of the tree and labeled the nodes by pie charts depicting the

194    distribution in chemical superclasses (Fig. S4a) and classes (Fig. S4b) of its tips. The molecular

195    fingerprint-based hierarchy of chemical features agreed well with ClassyFire taxonomy

196    assignment, further demonstrating that molecular fingerprints can meaningfully capture

197    structural relationships among molecules in a hierarchical manner. Furthermore, Qemistree

198    coupled the chemical tree to sample metadata, revealing distinct chemical classes expected for

199    each sample type. Branches representing acyl-carnitines were exclusively found in animal

200    products (shades of blue; Fig. 3a). In contrast, honey, although categorized as an animal product,

201    shared most of its chemical space with plant products, reflective of the plant nectar and pollen-

202    based diet of honey bees. We observed a clade of flavonoids in both plant products and honey

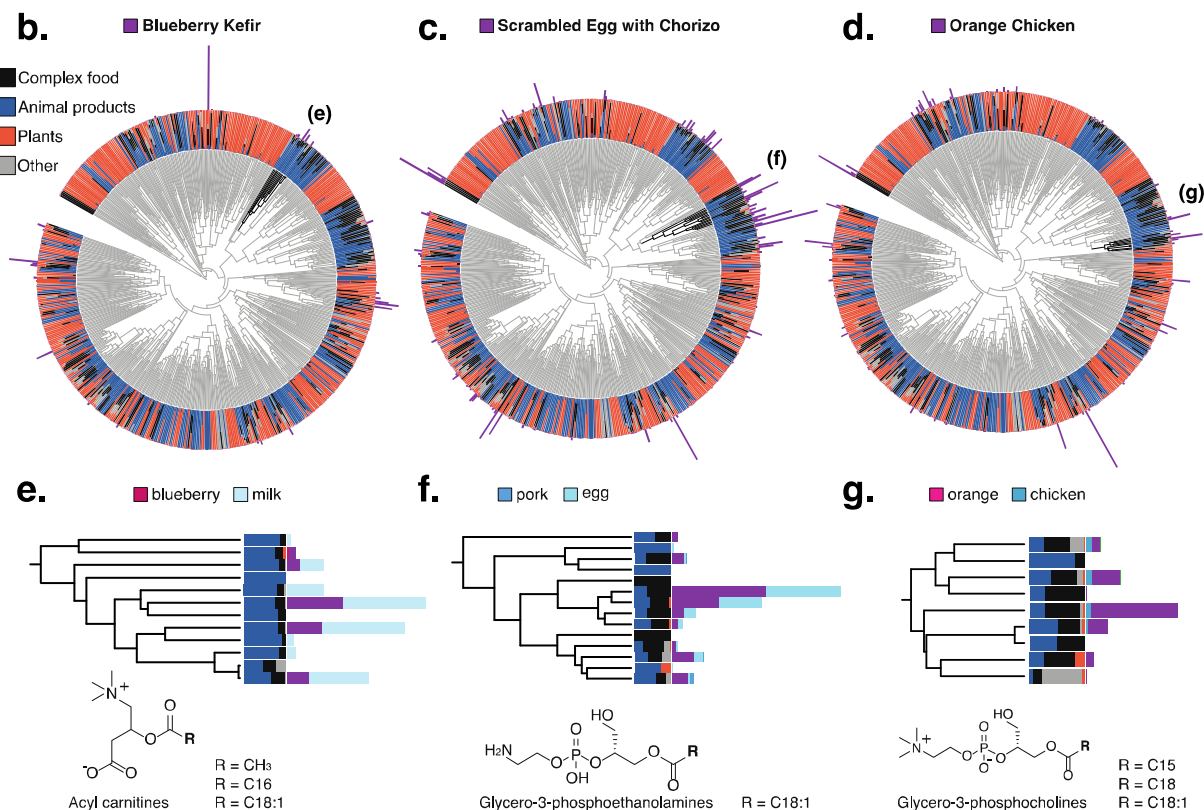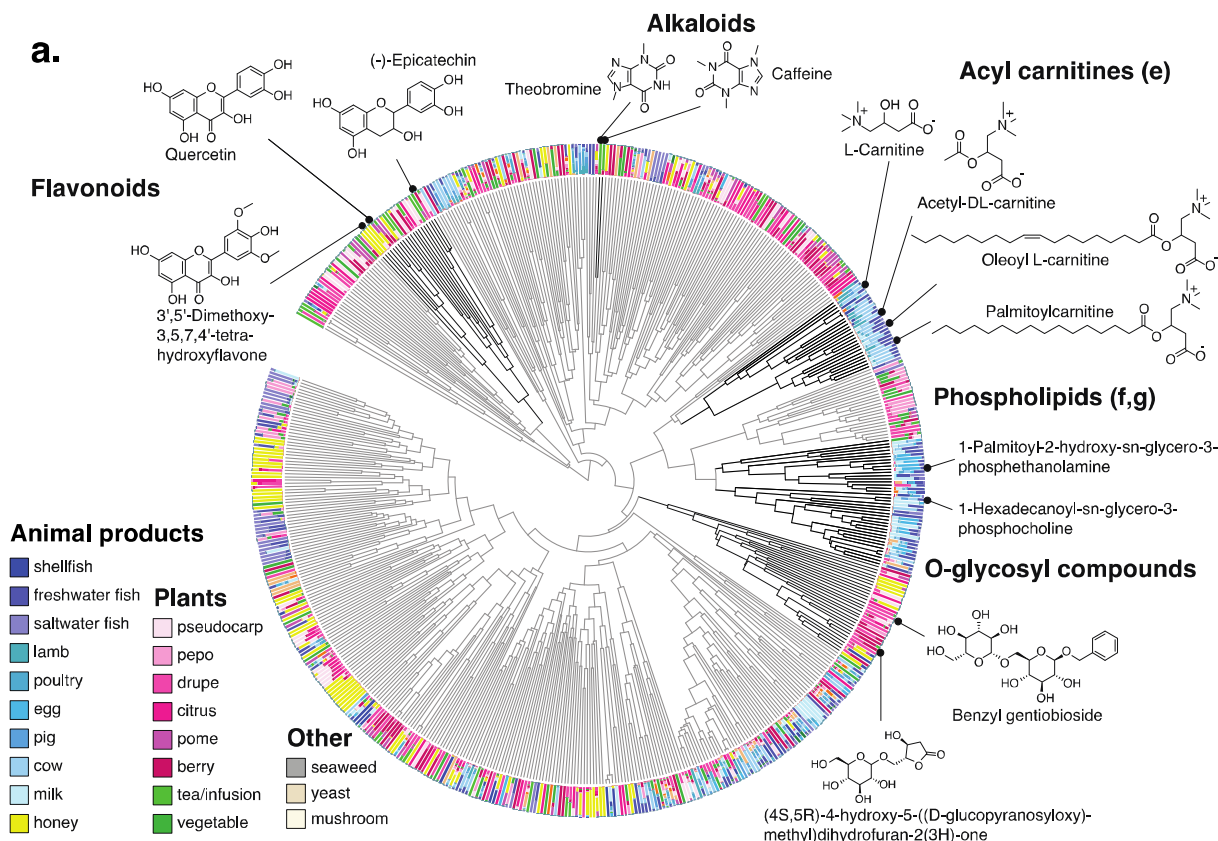203    (Figs. 3a, S4b), but no other animal-based foods.

205 **Figure 3: A chemical hierarchy of food-derived compounds based on predicted molecular fingerprints. a)** A
206 chemical tree based on molecular fingerprints representing the structural relationships between chemical features
207 (tree tips) detected in food products (single ingredient i.e. simple foods; N=119). The tree is pruned to only keep tips
208 that were assigned a structural annotation (SMILES) by either MS/MS spectral library match or *in silico* using
209 CSI:FingerID. All structures shown are spectral reference library matches obtained from feature-based molecular
210 networking in GNPS: (https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=ceb28a199d6b4f4fbf08490d9c96d631;
211 MSI level 3 annotation[27]). The outer ring shows the relative abundance of each compound across a diverse range of
212 food sources (panel a legend; parsed at 'sample_type_group4' of the Global FoodOmics Project ontology). We
213 highlight clusters of compounds that are characteristic of specific food sources. For example, theobromine and
214 caffeine are two closely related xanthine compounds (top center); they are primarily seen in teas (light green
215 samples) and coffee beans (berry; purple). Similarly, acyl-carnitines and phospholipids (top right) are unique to
216 different animal products (blues). We note that honey (highlighted in yellow), although annotated as an animal
217 product, contains compounds that are primarily seen in plant sources (flavonoids, O-glycosyl compounds) and no
218 other animal products. Flavonoids (top left) are observed in a range of fruit, vegetable, and honey samples (but no
219 other animal products). **(b-d)** A hierarchy of the compounds observed in simple foods (above) and seven complex
220 samples: two meals of orange chicken, a cooked cucumber and the sauce from a meal (schmorgurken), sour cream,
221 blueberry kefir, and egg scramble with chorizo (N=126). The inner ring shows the relative abundance of each
222 compound across simple animal products, plant products, fungi and algae (other) and the 7 complex foods (black).
223 The absolute abundances of compounds in blueberry kefir (b), scrambled eggs with chorizo (c), and orange chicken
224 (d) (outer bars) are overlaid on the tree to illustrate the shared and unique chemistry of complex foods. A compound
225 subtree characteristic of each complex food in the tree is highlighted (black) and zoomed in **(e-g)**. (e) A subtree
226 showing the absolute abundance of acyl carnitines in blueberry kefir and its primary ingredients (blueberry and
227 milk). Similar subtrees showing phosphoethanolamine in scrambled eggs with chorizo (f), and phosphocholine in
228 orange chicken (g)**.**
229

230 While it is expected that a complex food such as blueberry kefir contains molecules from both

231 blueberries and dairy, we can now visualize how individual ingredients and food preparation

232 contribute to the chemical composition of complex foods. We noted that metabolite signatures

233 that stem directly from particular ingredients, such as phosphoethanolamine from eggs, are

234 present in egg scramble (Fig. 3c), but not in the other two foods highlighted (Fig. 3b and d). We

235 can also observe the addition of ingredients in foods that were not listed as present in the initial

236 set of ingredients. We were able to retrieve that there is black pepper in the egg scramble with

237 chorizo and orange chicken, but that this signal is absent from the blueberry kefir (Fig. S5).

238

239 We show that our tree-based approach coherently captures chemical ontologies and relationships

240 among molecules and samples in various publicly available datasets. Qemistree depends on

241 representing chemical features as molecular fingerprints, and shares limitations with the

242 underlying fingerprint prediction tool CSI:FingerID. For example, fingerprint prediction depends

243 on the quality and coverage of MS/MS spectral databases available for training the predictive

244 models, and these will improve as databases are enriched with more compound classes.

245 Qemistree is also applicable in negative ionization mode; however, less molecular fingerprints

246    can be confidently predicted due to less publicly available reference spectra, resulting in less

247    extensive trees.

248

249    In summary, we introduce a new tree-based approach for computing and representing chemical

250    features detected in untargeted metabolomics studies. A hierarchy enables us to leverage existing

251    tree-based tools, and can be augmented with structural and environmental annotations, greatly

252    facilitating analysis and interpretation. We anticipate that Qemistree, as a data organization

253    strategy, will be broadly applicable across fields that perform global chemical analysis, from

254    medicine to environmental microbiology to food science, and well beyond the examples shown

255    here.

256

257 Data availability

258 The mass spectrometry data, metadata, and methods for the evaluation dataset have been

259 deposited on the GNPS/MassIVE public repository[2,33] under the accession number

260 MSV000083306. The parameters used for molecular networking are available on GNPS:

261 https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=efda476c72724b29a91693a108fa5a9d. The

262 chemical hierarchy generated by Qemistree (version 2020.1.2) is available on iTOL[24]:

263 https://itol.embl.de/tree/709513416494381587432576.

264 The mass spectrometry data, metadata, and methods for Global Foodomics dataset have been

265 deposited on the GNPS/MassIVE public repository[2,33] under the accession number

266 MSV000085226. The parameters used for molecular networking are available on GNPS:

267 https://gnps.ucsd.edu/ProteoSAFe/status.jsp?task=ceb28a199d6b4f4fbf08490d9c96d631. The

268 chemical hierarchy generated by Qemistree (version 2020.1.2) is available on iTOL[24]:

269 https://itol.embl.de/tree/13711034118313741584046018.

270 Code availability

271 All source code is publicly available under BSD-2-Clause on GitHub:

272 https://github.com/biocore/q2-qemistree. Qemistree is also available as an advanced analysis

273 workflow on GNPS: https://ccms-ucsd.github.io/GNPSDocumentation/qemistree/

274

275 Acknowledgments

283 Conflict of Interests

284 Mingxun Wang is a founder of Ometa Labs LLC.

285    Pieter C. Dorrestein is a scientific advisor for Sirenas LLC.

286    Kai Dührkop, Marcus Ludwig, Markus Fleischauer and Sebastian Böcker are founders of Bright

287    Giant GmbH.

288

289    References

290    1.   Watrous, J. *et al.* Mass spectral molecular networking of living microbial colonies. *Proc.*

291         *Natl. Acad. Sci. U. S. A.* **109**, E1743–52 (2012).

292    2.   Wang, M. *et al.* Sharing and community curation of mass spectrometry data with Global

293         Natural Products Social Molecular Networking. *Nat. Biotechnol.* **34**, 828–837 (2016).

294    3.   Fox Ramos, A. E., Evanno, L., Poupon, E., Champy, P. & Beniddir, M. A. Natural products

295         targeting strategies involving molecular networking: different manners, one goal. *Nat. Prod.*

296         *Rep.* **36**, 960–980 (2019).

297    4.   Böcker, S. & Dührkop, K. Fragmentation trees reloaded. *J. Cheminform.* **8**, 5 (2016).

298    5.   Rasche, F. *et al.* Identifying the unknowns by aligning fragmentation trees. *Anal. Chem.* **84**,

299         3417–3426 (2012).

300    6.   Washburne, A. D. *et al.* Phylogenetic factorization of compositional data yields lineage-

301         level associations in microbiome datasets. *PeerJ* **5**, e2969 (2017).

302    7.   Faith, D. P. Conservation evaluation and phylogenetic diversity. *Biological Conservation*

303         vol. 61 1–10 (1992).

304    8.   Janssen, S. *et al.* Phylogenetic Placement of Exact Amplicon Sequences Improves

305         Associations with Clinical Information. *mSystems* **3**, (2018).

306    9.   McDonald, D. *et al.* Striped UniFrac: enabling microbiome analysis at unprecedented scale.

307         *Nat. Methods* **15**, 847–848 (2018).

308    10.  Willett, P. Similarity-based virtual screening using 2D fingerprints. *Drug Discov. Today* **11**,

309     1046–1053 (2006).

310   11. Heinonen, M., Shen, H., Zamboni, N. & Rousu, J. Metabolite identification and molecular

311       fingerprint prediction through machine learning. *Bioinformatics* **28**, 2333–2341 (2012).

312   12. Laponogov, I., Sadawi, N., Galea, D., Mirnezami, R. & Veselkov, K. A. ChemDistiller: an

313       engine for metabolite annotation in mass spectrometry. *Bioinformatics* vol. 34 2096–2102

314       (2018).

315   13. Dührkop, K., Shen, H., Meusel, M., Rousu, J. & Böcker, S. Searching molecular structure

316       databases with tandem mass spectra using CSI:FingerID. *Proc. Natl. Acad. Sci. U. S. A.*

317       **112**, 12580–12585 (2015).

318   14. Fan, Z., Ghaffari, K., Alley, A. & Ressom, H. W. Metabolite Identification Using Artificial

319       Neural Network. *2019 IEEE International Conference on Bioinformatics and Biomedicine*

320       *(BIBM)* (2019) doi:10.1109/bibm47256.2019.8983190.

321   15. Li, Y., Kuhn, M., Gavin, A.-C. & Bork, P. Identification of metabolites from tandem mass

322       spectra with a machine learning approach utilizing structural features. *Bioinformatics* **36**,

323       1213–1218 (2020).

324   16. Dührkop, K. *et al.* SIRIUS 4: a rapid tool for turning tandem mass spectra into metabolite

325       structure information. *Nat. Methods* **16**, 299–302 (2019).

326   17. Pluskal, T., Castillo, S., Villar-Briones, A. & Oresic, M. MZmine 2: modular framework for

327       processing, visualizing, and analyzing mass spectrometry-based molecular profile data.

328       *BMC Bioinformatics* **11**, 395 (2010).

329   18. Nothias, L. F. *et al.* Feature-based Molecular Networking in the GNPS Analysis

330       Environment. *bioRxiv* 812404 (2019) doi:10.1101/812404.

331   19. Treutler, H. *et al.* Discovering Regulated Metabolite Families in Untargeted Metabolomics

332   Studies. *Anal. Chem.* **88**, 8082–8090 (2016).

333 20. Depke, T., Franke, R. & Brönstrup, M. Clustering of MS2 spectra using unsupervised

334   methods to aid the identification of secondary metabolites from Pseudomonas aeruginosa.

335   *Journal of Chromatography B* vol. 1071 19–28 (2017).

336 21. Rawlinson, C. *et al.* Hierarchical clustering of MS/MS spectra from the firefly metabolome

337   identifies new lucibufagin compounds. *Sci. Rep.* **10**, 6043 (2020).

338 22. Schymanski, E. L. *et al.* Critical Assessment of Small Molecule Identification 2016:

339   automated methods. *J. Cheminform.* **9**, 22 (2017).

340 23. Feunang, Y. D. *et al.* ClassyFire: automated chemical classification with a comprehensive,

341   computable taxonomy. *J. Cheminform.* **8**, 1–20 (2016).

342 24. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v4: recent updates and new

343   developments. *Nucleic Acids Res.* **47**, W256–W259 (2019).

344 25. Bolyen, E. *et al.* Reproducible, interactive, scalable and extensible microbiome data science

345   using QIIME 2. *Nat. Biotechnol.* **37**, 852–857 (2019).

346 26. Morton, J. T. *et al.* Learning representations of microbe-metabolite interactions. *Nat.*

347   *Methods* **16**, 1306–1314 (2019).

348 27. Sumner, L. W. *et al.* Proposed minimum reporting standards for chemical analysis.

349   *Metabolomics* vol. 3 211–221 (2007).

350 28. Sedio, B. E., Rojas Echeverri, J. C., Boya P., C. A. & Joseph Wright, S. Sources of

351   variation in foliar secondary chemistry in a tropical forest tree community. *Ecology* vol. 98

352   616–623 (2017).

353 29. Bray, J. R., Roger Bray, J. & Curtis, J. T. An Ordination of the Upland Forest Communities

354   of Southern Wisconsin. *Ecological Monographs* vol. 27 325–349 (1957).

355    30.  Gauglitz, J. M. *et al.* Metabolome-informed microbiome analysis refines metadata

356          classifications and reveals unexpected medication transfer in captive cheetahs. *bioRxiv*

357          790063 (2019) doi:10.1101/790063.

358    31.  Thompson, L. R. *et al.* A communal catalogue reveals Earth's multiscale microbial

359          diversity. *Nature* **551**, 457–463 (2017).

360    32.  Morton, J. T. *et al.* Establishing microbial composition measurement standards with

361          reference frames. *Nat. Commun.* **10**, 2719 (2019).

362    33.  Wang, M. *et al.* Assembling the Community-Scale Discoverable Human Proteome. *Cell*

363          *Syst* **7**, 412–421.e5 (2018).

364    34.  Ludwig, M. *et al.* ZODIAC: database-independent molecular formula annotation using

365          Gibbs sampling reveals unknown small molecules. *bioRxiv* 842740 (2019)

366          doi:10.1101/842740.

367    35.  Simón-Manso, Y. *et al.* Metabolite profiling of a NIST Standard Reference Material for

368          human plasma (SRM 1950): GC-MS, LC-MS, NMR, and clinical laboratory analyses,

369          libraries, and web-based resources. *Anal. Chem.* **85**, 11725–11731 (2013).

370    36.  McDonald, D. *et al.* American Gut: an Open Platform for Citizen Science Microbiome

371          Research. *mSystems* **3**, (2018).

372    37.  Martens, L. *et al.* mzML--a community standard for mass spectrometry data. *Mol. Cell.*

373          *Proteomics* **10**, R110.000133 (2011).

374    38.  Chambers, M. C. *et al.* A cross-platform toolkit for mass spectrometry and proteomics. *Nat.*

375          *Biotechnol.* **30**, 918–920 (2012).

376