

1 **Title:** Reproducible Molecular Networking Of Untargeted Mass Spectrometry Data Using  
2 GNPS.

3  
4 **Authors:** Allegra T. Aron<sup>1</sup> #, Emily C. Gentry<sup>1</sup> #, Kerry L. McPhail<sup>2</sup> #, Louis Felix Nothias<sup>1</sup>,  
5 Mélissa Nothias-Esposito<sup>1</sup>, Amina Bouslimani<sup>1</sup>, Daniel Petras<sup>1,5</sup>, Julia M. Gauglitz<sup>1</sup>, Nicole  
6 Sikora<sup>1</sup>, Fernando Vargas<sup>1</sup>, Justin J. J. van der Hooft<sup>3</sup>, Madeleine Ernst<sup>1</sup>, Kyo Bin Kang<sup>4</sup>,  
7 Christine M. Aceves<sup>1</sup>, Andrés Mauricio Caraballo-Rodríguez<sup>1</sup>, Irina Koester<sup>1,5</sup>, Kelly C.  
8 Weldon<sup>1</sup>, Samuel Bertrand<sup>6,7</sup>, Catherine Roullier<sup>6</sup>, Kunyang Sun<sup>1</sup>, Richard M. Tehan<sup>2</sup>,  
9 Christopher A. Boya<sup>8,9</sup>, Christian Martin H.<sup>8</sup>, Marcelino Gutiérrez<sup>8</sup>, Aldo Moreno Ulloa<sup>10</sup>,  
10 Javier Andres Tejeda Mora<sup>10</sup>, Randy Mojica-Flores<sup>8,11</sup>, Johant Lakey-Beitia<sup>8</sup>, Victor  
11 Vásquez-Chaves<sup>12</sup>, Yilue Zhang<sup>13</sup>, Angela I. Calderon<sup>13</sup>, Nicole Tayler<sup>8,9</sup>, Robert A.  
12 Keyzers<sup>14</sup>, Fidele Tugizimana<sup>15</sup>, Nombuso Ndlovu<sup>15</sup>, Alexander A. Aksenov<sup>1</sup>, Alan  
13 Jarmusch<sup>1</sup>, Robin Schmid<sup>16</sup>, Andrew W. Truman<sup>17</sup>, Nuno Bandeira<sup>18\*</sup>, Mingxun Wang<sup>1\*</sup>,  
14 Pieter C Dorrestein<sup>1, 19-21\*</sup>

15  
16 **Affiliations:** <sup>1</sup>Skaggs School of Pharmacy and Pharmaceutical Sciences, University of  
17 California San Diego, La Jolla, California, USA. <sup>2</sup>Department of Pharmaceutical Sciences,  
18 College of Pharmacy, Oregon State University, Corvallis, Oregon, USA. <sup>3</sup>Bioinformatics  
19 Group, Wageningen University, Wageningen 6708 PB, The Netherlands. <sup>4</sup>College of  
20 Pharmacy, Sookmyung Women's University, Cheongpa-ro 47-gil 100, Yongsan-gu, Seoul  
21 04310, Korea. <sup>5</sup>Scripps Institution of Oceanography, University of California San Diego,  
22 La Jolla, California, USA. <sup>6</sup>Groupe Mer, Molécules, Santé-EA 2160, UFR des Sciences  
23 Pharmaceutiques et Biologiques, Université de Nantes, 44035 Nantes, France.  
24 <sup>7</sup>ThalassOMICS Metabolomics Facility, Plateforme Corsaire, Biogenouest, 44035 Nantes,  
25 France. <sup>8</sup>Centro de Biodiversidad y Descubrimiento de Drogas, Instituto de  
26 Investigaciones Científicas y Servicios de Alta Tecnología (INDICASAT AIP), Panamá,  
27 Apartado 0843-01103, República de Panamá. <sup>9</sup>Department of Biotechnology, Acharya  
28 Nagarjuna University, Guntur, Nagarjuna Nagar-522 510, India. <sup>10</sup>Biomedical Innovation  
29 Department, CICESE, México. <sup>11</sup>Universidad Autónoma de Chiriquí (UNACHI), Mexico.  
30 <sup>12</sup>Centro de Investigaciones en Productos Naturales (CIPRONA), Universidad de Costa  
31 Rica, San José, Costa Rica. <sup>13</sup>Harrison School of Pharmacy, Auburn University, Auburn,  
32 Alabama, USA. <sup>14</sup>School of Chemical & Physical Sciences, Victoria University of  
33 Wellington, Wellington, New Zealand. <sup>15</sup>Centre for Plant Metabolomics Research,  
34 Department of Biochemistry, University of Johannesburg, Auckland Park 2006, South  
35 Africa. <sup>16</sup>Institute of Inorganic and Analytical Chemistry, University of Münster, 48149  
36 Münster, Germany. <sup>17</sup>Department of Molecular Microbiology, John Innes Centre, Norwich,  
37 NR4 7UH, U.K. <sup>18</sup>Computer Science and Engineering, University of California San Diego,  
38 La Jolla, California, USA. <sup>19</sup>Center for Computational Mass Spectrometry, University of  
39 California San Diego, La Jolla, California, USA. <sup>20</sup>Department of Pharmacology, University  
40 of California San Diego, La Jolla, California, USA. <sup>21</sup>Department of Pediatrics, University  
41 of California San Diego, La Jolla, California, USA. #These authors contributed equally to  
42 this work. Correspondence should be addressed to N.B.(nbandeira@ucsd.edu, M.W.  
43 (miw023@ucsd.edu) or P.C.D. (pdorrestein@ucsd.edu)

44  
45 **Author contributions:** Design and oversight of the project: P.C.D., M.W., N.B. Instrument  
46 acquisition parameters: A.T.A., E.C.G., K.L.M., R.M.T., K.B.K., S.B., C.R., A.W.T., F.T.,  
47 N.N., A.M.U. Data conversion and upload: K.L.M., E.C.G., A.T.A., J.J.J. v.d.H., M.E. GNPS  
48 documentation: M.W., L.F.N., E.C.G., A.T.A., K.L.M., J.J.J.v.d.H., M.E, M.N.-E. Cytoscape

49 documentation: M.N-E., F.V., I.K., A.M.C-R. Metadata curation: J.M.G., C.M.A., F.V.,  
50 A.M.C-R. Mass spectra annotations: D.P, R.S., M.E. Theoretical tools and advanced  
51 features, statistical analysis: L.F.N., A.A. Supplementary information: A.T.A., N.S., E.C.G.,  
52 K.L.M., M.E. Testing the workflows described and improving the descriptions: A.I.C,  
53 A.M.U, J.A.T.M, C.M.H., C.A.B.P., M.G., V.V-C., J.L-B., R.M-F., M.E.

54

55 Authors names and emails:

56 Allegra T. Aron ([alaron@ucsd.edu](mailto:alaron@ucsd.edu))  
57 Emily C. Gentry ([emgentry@ucsd.edu](mailto:emgentry@ucsd.edu))  
58 Kerry L. McPhail ([kerry.mcphail@oregonstate.edu](mailto:kerry.mcphail@oregonstate.edu))  
59 Louis-Felix Nothias ([lnothiasscaglia@ucsd.edu](mailto:lnothiasscaglia@ucsd.edu))  
60 Julia M. Gauglitz ([jgauglitz@ucsd.edu](mailto:jgauglitz@ucsd.edu))  
61 Christine M. Aceves ([caceves@ucsd.edu](mailto:caceves@ucsd.edu))  
62 Fernando Vargas ([fernando.vargas0341@gmail.com](mailto:fernando.vargas0341@gmail.com))  
63 Amina Bouslimani ([abouslimani@ucsd.edu](mailto:abouslimani@ucsd.edu))  
64 Justin J. J. van der Hoof ([justin.vanderhoof@wur.nl](mailto:justin.vanderhoof@wur.nl))  
65 Kyo Bin Kang ([kbkang@sookmyung.ac.kr](mailto:kbkang@sookmyung.ac.kr))  
66 Andrés Mauricio Caraballo-Rodríguez ([acaraballo-rodriguez@ucsd.edu](mailto:acaraballo-rodriguez@ucsd.edu))  
67 Irina Koester ([ikoester@ucsd.edu](mailto:ikoester@ucsd.edu))  
68 Kelly C. Weldon ([kcweldon@ucsd.edu](mailto:kcweldon@ucsd.edu))  
69 Daniel Petras ([dpetras@ucsd.edu](mailto:dpetras@ucsd.edu))  
70 Samuel Bertrand ([Samuel.Bertrand@univ-nantes.fr](mailto:Samuel.Bertrand@univ-nantes.fr))  
71 Catherine Roullier ([Catherine.Roullier@univ-nantes.fr](mailto:Catherine.Roullier@univ-nantes.fr))  
72 Madeleine Ernst ([mernst@ucsd.edu](mailto:mernst@ucsd.edu))  
73 Kunyang Sun ([ksun@ucsd.edu](mailto:ksun@ucsd.edu))  
74 Richard M. Tehan ([tehanr@oregonstate.edu](mailto:tehanr@oregonstate.edu))  
75 Cristopher A. Boya P. ([c.boya@indicat.org.pa](mailto:c.boya@indicat.org.pa))  
76 Christian Martin H. ([cmartin@indicat.org.pa](mailto:cmartin@indicat.org.pa))  
77 Marcelino Gutiérrez ([mgutierrez@indicat.org.pa](mailto:mgutierrez@indicat.org.pa))  
78 Aldo Moreno Ulloa ([amoreno@cicese.mx](mailto:amoreno@cicese.mx))  
79 Javier Andres Tejeda Mora ([andres.android@gmail.com](mailto:andres.android@gmail.com))  
80 Randy Mojica-Flores ([wendel2506@gmail.com](mailto:wendel2506@gmail.com))  
81 Johant Lakey-Beitia ([jlakey@indicat.prg.pa](mailto:jlakey@indicat.prg.pa))  
82 Victor Vásquez-Chaves ([vvasquezch@gmail.com](mailto:vvasquezch@gmail.com))  
83 Angela I. Calderon ([aic0001@auburn.edu](mailto:aic0001@auburn.edu))  
84 Nicole Tayler ([ntayler@indicat.org.pa](mailto:ntayler@indicat.org.pa))  
85 Robert A. Keyzers ([robert.Keyzers@vuw.ac.nz](mailto:robert.Keyzers@vuw.ac.nz))  
86 Fidele Tugizimana ([ftugizimana@uj.ac.za](mailto:ftugizimana@uj.ac.za))  
87 Nombuso Ndlovu ([nndlovu@uj.ac.za](mailto:nndlovu@uj.ac.za))  
88 Nicole Sikora ([nsikora@ucsd.edu](mailto:nsikora@ucsd.edu))  
89 Alexander Aksenov ([aaaksenov@ucsd.edu](mailto:aaaksenov@ucsd.edu))  
90 Alan Jarmusch ([ajarmusch@ucsd.edu](mailto:ajarmusch@ucsd.edu))  
91 Robin Schmid ([robinschmid@uni-muenster.de](mailto:robinschmid@uni-muenster.de))  
92 Andrew W. Truman ([Andrew.Truman@jic.ac.uk](mailto:Andrew.Truman@jic.ac.uk))  
93 Nuno Bandeira ([bandeira@ucsd.edu](mailto:bandeira@ucsd.edu))  
94 Mingxun Wang ([miw023@ucsd.edu](mailto:miw023@ucsd.edu))  
95 Pieter C Dorrestein ([pdorrestein@ucsd.edu](mailto:pdorrestein@ucsd.edu))

96

97 **Abstract:** Global Natural Product Social (GNPS) Molecular Networking is an interactive  
98 online chemistry-focused mass spectrometry data curation and analysis infrastructure.  
99 The goal of GNPS is to provide as much chemical insight for an untargeted tandem mass  
100 spectrometry data set as possible and to connect this chemical insight to the underlying  
101 biological questions a user wishes to address. This can be performed within one  
102 experiment or at the repository scale. GNPS not only serves as a public data repository  
103 for untargeted tandem mass spectrometry data with the sample information (metadata), it  
104 also captures community knowledge that is disseminated *via* living data across all public  
105 data. One of the main analysis tools used by the GNPS community is molecular  
106 networking. Molecular networking creates a structured data table that reflects the chemical  
107 space from tandem mass spectrometry experiments *via* computing the relationships of the  
108 tandem mass spectra through spectral similarity. This protocol provides step-by-step  
109 instructions for creating reproducible high-quality molecular networks. For training  
110 purposes, the reader is led through the protocol from recalling a public dataset and its  
111 sample information to creating and interpreting a molecular network. Each data analysis  
112 job can be shared or cloned to disseminate the knowledge gained, thus propagating  
113 information that can lead to the discovery of molecules, metabolic pathways and  
114 ecosystem/community interactions.

115

116 **1.0 Introduction:** Molecular networking for the analysis of tandem mass spectra of small  
117 molecules was introduced in 2012<sup>1</sup>. Upon its introduction, molecular networking was  
118 compared to sequencing of environmental DNA to study the microbial communities  
119 present in diverse ecosystems<sup>2</sup>. For the first time we were able to get a map of the  
120 chemical diversity that is observed in an untargeted mass spectrometry experiment. In  
121 addition to providing unprecedented systems-level views of the chemical space in various  
122 environments, molecular networking has aided structure elucidation of many compounds<sup>3-</sup>  
123 <sup>9</sup>.

124 The foundation of molecular networking is pairwise spectral alignment using a  
125 modified cosine spectral similarity algorithm originally intended to discover modified forms  
126 of peptides and proteins<sup>10</sup>. In a modified spectral similarity search, not only are  
127 fragmentation spectra (MS<sup>2</sup>) from ions at identical *m/z* compared, but also MS<sup>2</sup> spectra  
128 that are offset by the same *m/z* difference as the precursor ion. By eliminating the amino  
129 acid filtering from the original spectral alignment algorithms, it became possible to extend  
130 spectral similarity to any set of MS<sup>2</sup> spectra, including those from small molecules and  
131 natural products. When a pairwise spectral similarity search/alignment is performed, each  
132 MS<sup>2</sup> spectrum in a given dataset is compared against every other, and a network of MS<sup>2</sup>  
133 spectral relations is obtained, from which molecular networks are created (**Fig. 1**).  
134 Molecular networking build on the fundamental observation that two structurally related  
135 molecules share fragment ion patterns when subjected to MS<sup>2</sup> fragmentation methods  
136 such as collision induced dissociation (CID). In order to make the molecular networking  
137 algorithm accessible to the scientific community, its script was converted to a web-based  
138 platform backed by a supercomputer. This enabled the creation of a community  
139 infrastructure supporting both a database and knowledge-base around the needs of the  
140 community. The result was the Global Natural Products Social (GNPS) Molecular  
141 Networking community effort that started in 2014 and was published in 2016. The user  
142 base has expanded to 49 of 50 states in the United States and worldwide to over 150  
143 countries<sup>11</sup>. GNPS is currently widely used by scientists working in industry, academia and  
144 government in the fields of biomedical research, environmental science, ecology,

145 forensics, microbiology, chemistry, and others. This crowdsourced, community-driven  
 146 analysis infrastructure not only facilitates data and knowledge storage but also enables  
 147 knowledge capture, sharing, dissemination and data driven social networking while  
 148 promoting reproducible data analysis. Moreover, GNPS can be accessed on a computer  
 149 or on any mobile device connected to the internet making any public data set readily  
 150 accessible for analysis. While there are many analysis tools available within the GNPS  
 151 infrastructure, molecular networking is the most frequently used tool. Other tools available  
 152 on GNPS such as network annotation propagation (NAP) briefly discussed in section 3.5.

153

154 To create a molecular network, GNPS first aligns each MS<sup>2</sup> spectrum in a dataset to each  
 155 of the others, and assigns a *cosine score* to each combination to describe their similarity  
 156 (**Fig. 1**). Identical mass are collapsed based on a hierarchical cosine clustering algorithm  
 157 into a single *node* or *consensus cluster* due to the high similarity of their fragment ions.  
 158 This is accomplished using the MS-Cluster algorithm<sup>12</sup>. Structurally related molecules yield  
 159 comparable MS<sup>2</sup> spectra due to commonalities in their gas phase chemistry<sup>13</sup>, and are  
 160 represented by separate nodes that connect within the network via *edges*. Each  
 161 consensus spectrum (node) is then queried against spectral library databases to assign  
 162 putative known molecules within a network.

163

164

**Table 1.** Terminology

Term	Definition
<i>annotation</i>	The process of attributing a putative chemical structure to a detected molecule. The level of annotations from spectral matches are considered level 2 or 3 according to the 2007 Metabolomics Standards Initiative <sup>14</sup> .
<i>bucket table</i>	A tab separated table (.tsv file format) downloadable from the GNPS interface, which shows per sample summed precursor ion intensities per MS <sup>2</sup> ion. Pie charts generated in visualization tools are based off of intensities in the bucket table.
<i>cluster index</i>	Reference identification number for a MS <sup>2</sup> consensus cluster. In Cytoscape this identification number is also called 'shared name'.
<i>consensus cluster</i>	A grouping of MS <sup>2</sup> spectra that are considered identical based on the MS-Cluster algorithm <sup>10,12</sup> . Since GNPS brings together approaches from different scientific communities, there are terms such as "cluster" that have different meanings. Thus, the context in which the term is used should be considered. The term 'consensus cluster' refers to the grouping of MS <sup>2</sup> spectra into a node and is different from clusters of nodes in molecular networks as visualized in Cytoscape <sup>15,16</sup> .
<i>cosine score</i>	A value that represents the MS <sup>2</sup> spectral similarity between two nodes in the molecular network, where a cosine score of 1 represents identical spectra and a cosine score of 0 denotes no

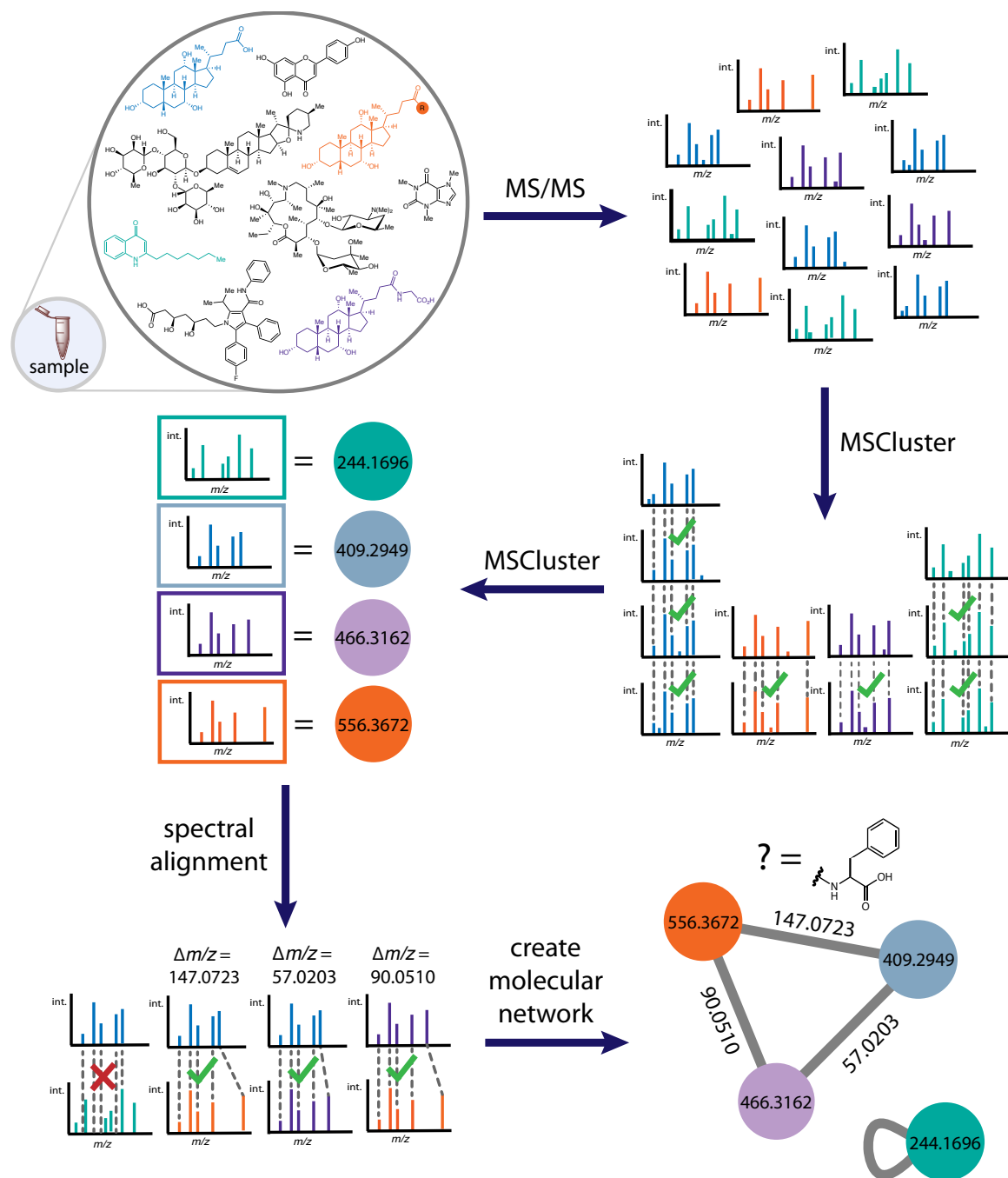
	similarity at all. The cosine score takes into account precursor ion, fragment ions as well as peak intensities <sup>1</sup> .
<i>DDA</i>	Abbreviation for data-dependent acquisition; a method for tandem mass spectrometry data collection where the most intense MS <sup>1</sup> ions are iteratively selected for MS <sup>2</sup> fragmentation <sup>17</sup> .
<i>dereplication</i>	Rapid identification of previously characterized (known) molecules <sup>18</sup> .
<i>edge</i>	A line connecting nodes that represents related but not identical MS <sup>2</sup> spectra based on a cosine similarity score.
<i>identification</i>	Validation of a molecular assignment using an authentic chemical standard analyzed under the same experimental conditions as the sample containing the unknown compound. Molecular identification requires matching at least one physical characteristic, e.g. retention time, exact <i>m/z</i> , and MS <sup>2</sup> fragmentation pattern <sup>14, 19</sup> .
<i>LC</i>	Abbreviation for liquid chromatography; a method used to separate molecules in a mixture using a liquid mobile phase.
<i>natural product</i>	A small molecule (< 2000 Da) produced by a biological source <sup>20</sup> .
<i>m/z</i>	Mass-to-charge ratio, a dimensionless quantity resulting from dividing the mass number of an ion by its charge number. <sup>21</sup>
<i>molecular network</i>	A map of all nodes illustrating connectivity that represents the chemical space detected in the experiment.
<i>molecular networking</i>	A computational approach that organizes MS <sup>2</sup> data based on spectral similarity, from which we can infer relationships in chemical structures <sup>1</sup> .
<i>MSCluster</i>	An algorithm used by GNPS to collapse nearly identical MS <sup>2</sup> spectra with the same precursor ion <i>m/z</i> into a single consensus spectrum.
<i>MS<sup>1</sup></i>	The collection of all precursor ions ( <i>m/z</i> ) and associated abundancies in a sample. MS <sup>1</sup> is the first stage of tandem mass spectrometry, where compounds can be further fragmented <sup>22, 23</sup> . See also <i>tandem MS</i> , <i>MS<sup>2</sup></i> , <i>MS/MS</i> .

<i>node</i>	A consensus cluster of identical MS <sup>2</sup> spectra that represent one molecule, or a single MS <sup>2</sup> spectrum if cluster size is 1.
<i>precursor ion (parent ion)</i>	The ionized form of a molecule that is selected for tandem MS fragmentation. In electrospray ionization, the parent ion is a synonym of precursor ion <sup>21</sup> .
<i>product ion (fragment ion)</i>	An ion originating from a gas-phase reaction of the precursor ion <sup>13</sup> .
<i>sample information (metadata)</i>	Data that provide basic information about the sample and descriptions to facilitate data analysis and interpretation. Examples of sample information include: the identification number, the source and origin of the sample collected, time, age, sex and date of collection.
<i>small molecule</i>	This protocol considers a molecule with a molecular weight < 1500 Da a small molecule.
<i>spectral alignment</i>	An algorithmic approach that aligns related spectra. This is the basis of molecular networking which relies on the assumption that two structurally related molecules share similarity in their MS <sup>2</sup> spectra <sup>1</sup> .
<i>spectral similarity</i>	The likeness of MS <sup>2</sup> spectra based on all or some of the following: precursor ion, fragment ions, and relative intensities of these peaks. Structurally related molecules tend to exhibit similar fragmentation <sup>13</sup> . In molecular networking spectral similarity is calculated through a modified cosine score.
<i>summed ion intensities</i>	Sum of precursor ion intensities in the MS <sup>2</sup> spectra for all ions with the same associated MS <sup>2</sup> detected by the mass spectrometer.
<i>tandem MS, MS/MS, MS<sup>2</sup></i>	Abbreviations for tandem mass spectrometry, which defines a technique where mass-selected ions are subjected to a second mass spectrometric analysis. In the first stage, also referred to as MS <sup>1</sup> , precursor ions are formed and detected. In the second stage, also referred to as MS <sup>2</sup> or MS/MS, precursor ions are fragmented resulting in a spectral fingerprint <sup>22, 23</sup> .

165

166 All mass spectrometry data used in GNPS, both in the private user workspace or data  
167 that are made public, is stored in MassIVE - an interactive virtual environment developed  
168 to facilitate and encourage the exchange of mass spectrometry data. MassIVE accepts  
169 data files (organized as datasets) and facilitates the sharing of datasets with a unique  
170 identifier; one can use this unique identifier as an accession number for publications. In  
171 addition, public datasets that the user publishes can, by choice of the depositor, have an  
172 associated DOI. Currently, MassIVE is an approved repository for the Journal of  
173 Proteome Research (<https://pubs.acs.org/journal/jprobs>) and Nature Partner Journals  
174 (<https://www.nature.com/sdata/policies/repositories#chem>) and is widely used as a  
175 repository for other journals<sup>24-33</sup>. GNPS-MassIVE has more than a thousand public

176 metabolomics datasets. The GNPS knowledge base includes 221,083 reference MS<sup>2</sup>  
177 spectra, provided by the GNPS community, spectral libraries generated for GNPS  
178 (GNPS-collections) and third party libraries<sup>11</sup>. Examples are LDB Lichen Database,  
179 MIADB Spectral Library, Sumner Spectral Library, CASMI Spectral Library, and  
180 Massbank, a large MS data library that is directly synced with GNPS. There are also tags  
181 and sample information (metadata) entries provided by the community in the GNPS  
182 knowledge base. Furthermore, all public data is periodically searched against the NIST  
183 2017 spectral library and high confidence spectral matches are annotated. GNPS-  
184 MASSIVE now performs more than 6,000 analysis jobs a month and more than 200,000  
185 page views (excluding developers), with the predominant analysis being molecular  
186 networking. As a result, GNPS based analysis has been used for the discovery of  
187 hundreds of new molecules in the last few years, ranging from immune regulators to  
188 antimicrobials, including antiviral agents and protease inhibitors<sup>9, 34-36</sup>. Here we provide a  
189 detailed protocol on how to generate a publishable and reproducible molecular network  
190 from a mass spectrometry dataset. This protocol will take the reader through the  
191 following steps: how to upload data, how to make the data public, how to subscribe to  
192 public data for living data updates, and how to reproducibly create publishable molecular  
193 networks using standardized sample information (metadata) through the GNPS  
194 infrastructure (**Fig. 1**).  
195

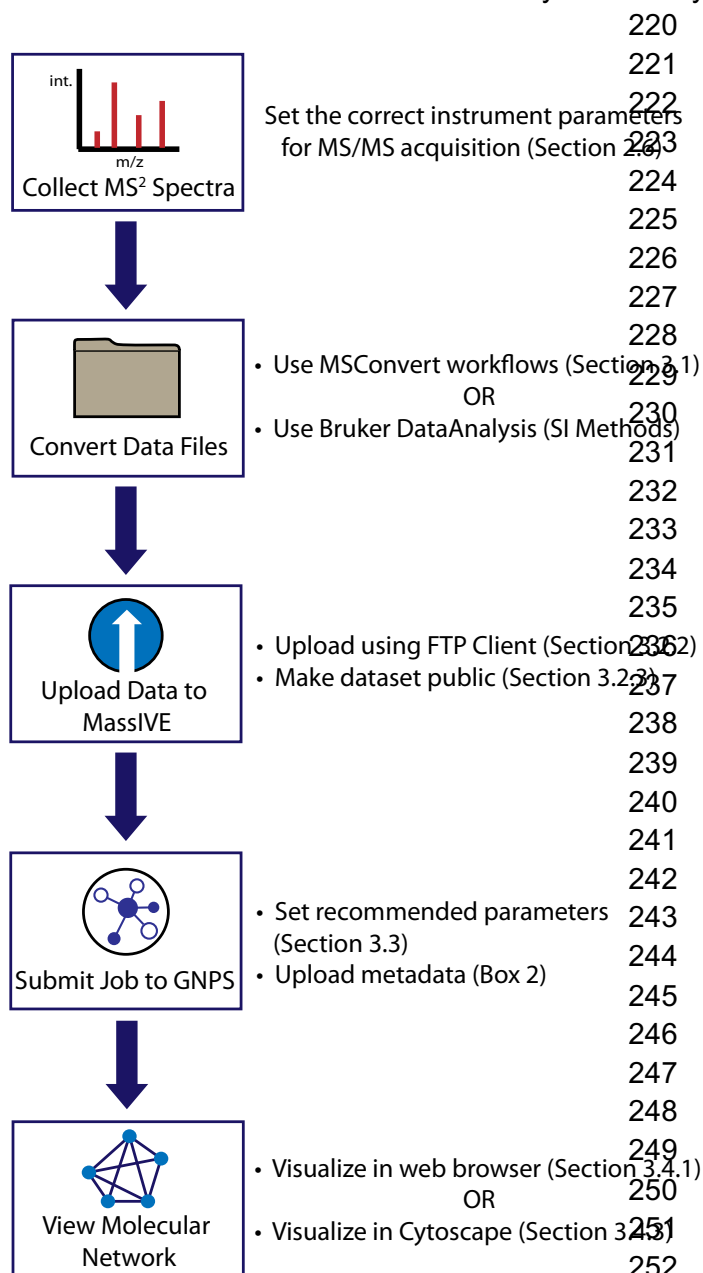


196  
 197  
 198  
 199  
 200  
 201  
 202  
 203  
 204  
 205  
 206  
 207

**Figure 1.** Schematic representation of the process for creating a molecular network from tandem mass spectra acquired for metabolites in complex sample mixtures. The colors are used to track how we go from molecules in a sample to nodes in the molecular network. We start by obtaining MS<sup>2</sup> spectra of all ionized molecules in the sample. MS-Cluster first aligns each MS<sup>2</sup> spectrum in a dataset to each of the others. Mass spectra from identical compounds coalesced using MS-Cluster<sup>12</sup> into a single *node* or *consensus cluster* due to the high similarity of their precursor ion and fragment ions. Subsequently a spectral alignment is performed enabling for similarity searches even when the precursor ion masses are not identical. This is accomplished using a modified cosine score where all the ions that differ by the mass difference of the two precursor ions are also considered.



208 Structurally-related molecules yield comparable MS<sup>2</sup> spectra due to commonalities in their  
 209 gas phase chemistry, and are represented by separate nodes that connect within the  
 210 network via *edges*. Each node is then queried against spectral libraries to assign putative  
 211 known molecules within a molecular network and unknowns can be propagated using  
 212 chemical rationale. For illustration purposes, the blue node with *m/z* 409.2949 is cholate,  
 213 *m/z* 446.3162 in purple is glycocholic acid (the user would discover this based on MS<sup>2</sup>  
 214 matches to a reference library) while the orange one is unknown but has a mass shift of  
 215 147.0723 Da. This is a typical mass shift of phenylalanine and thus a prediction can be  
 216 made that this is a phenylalanine conjugate of cholic acid. The difference between the  
 217 glycine and phenylalanine conjugate is 90.0510 Da and supports such structural  
 218 hypothesis. The self looped green node *m/z* 244.1696 is attributed to an unrelated  
 219 molecule and therefore does not have any structurally related molecule in the sample.



254 data acquisition, conversion, upload and networking to visualization. Readers following the  
 255 tutorial example can follow these steps to generate a publishable network.

256

257 Data collection and processing procedures will vary depending on the instrument  
258 available to the user. Although the user can modify any procedure to fit their specific goals,  
259 this protocol specifies a set of starting parameters for acquiring and converting data with  
260 various mass spectrometers, including AB Sciex, Agilent, Bruker, Shimadzu, Thermo  
261 Scientific, and Waters instruments. We also provide a protocol for the conversion of the  
262 data from each of these vendors to an open format (.mzXML, .mzML or MGF) that is usable  
263 within the GNPS-MassIVE infrastructure. Once the data is converted to the proper open  
264 format, the protocol describes how to upload data files to MassIVE, a public repository that  
265 enables community sharing of mass spectrometry data, using either a web browser or FTP  
266 client. The resulting datasets can then be subsequently submitted to GNPS for molecular  
267 networking analysis, wherein MS<sup>2</sup> spectra are organized in a network according to  
268 similarity and compared against a reference database to identify putative known molecules  
269 and 'molecular families' in the samples. Finally, visualization and analysis of GNPS-  
270 generated molecular networks can be performed either in the web browser itself or in  
271 [Cytoscape](#), an open-source software for visualizing complex networks<sup>37</sup>.

272

### 273 **1.2 Applications of the method**

274 GNPS molecular networking provides the ability to analyze and compare MS<sup>2</sup> spectra in  
275 one or more datasets acquired within the scope of a specific study, across datasets from  
276 multiple studies, and also to compare those datasets to all publicly available GNPS-  
277 MassIVE datasets, including community curated spectral libraries. In addition, ongoing  
278 contributions to spectral libraries and submissions of new public datasets enable  
279 continuous identification: the periodic and automated reanalysis of all public datasets.  
280 GNPS is being used to network data acquired on a number of different mass  
281 spectrometers in a wide variety of exploratory studies, with samples originating from  
282 diverse environments and used for varying purposes. These range from the indoor  
283 environment<sup>38-40</sup> to dissolved organic matter in the oceans<sup>41</sup>, from microbes in culture<sup>9, 42-  
284 45</sup> to mouse<sup>46</sup> or human microbiomes<sup>47, 48</sup> or infections<sup>49-51</sup>, from clinical samples<sup>32, 52, 53</sup> to  
285 plants<sup>54</sup>, algae<sup>55</sup>, sponges<sup>5, 56</sup> and corals<sup>57</sup>, as well as a number of other sample types<sup>26,  
286 58</sup>. Additionally, molecular networking has been applied to natural products discovery from  
287 a variety of organisms<sup>59-62</sup>, forensics<sup>63</sup>, small molecule identification<sup>64</sup> and biological  
288 discovery in hypothesis-driven research<sup>65</sup>. Furthermore, GNPS facilitates large-scale  
289 meta-analyses that can compare and potentially link studies from different laboratories by  
290 enabling rapid comparisons across multiple public datasets. Finally, to promote data  
291 analysis reproducibility, all analysis jobs are saved together with their parameters, which  
292 can be shared or cloned for reanalysis; no other platform provides this service.

293

### 294 **1.3 Alternative methods to this protocol**

295 Several aspects of the GNPS-based molecular networking protocol are provided  
296 elsewhere, but not previously as a coherent workflow in one package. There are several  
297 repositories where metabolomics data can be uploaded<sup>66-68</sup>. According to the OMICS  
298 discovery index, the most widely used are GNPS-MassIVE, Metabolomics workbench<sup>12</sup>  
299 and MetaboLights<sup>69, 70</sup>.

300

301 Mass spectral library searching, or comparing MS<sup>2</sup> spectra of compounds in a  
302 sample to reference data in order to annotate metabolites<sup>71</sup>, has been implemented  
303 extensively, and successfully, for decades. Numerous commercial and non-commercial  
MS<sup>2</sup> reference databases exist, such as the NIST/EPA/NIH Mass Spectral Library<sup>72</sup>,

304 METLIN<sup>73</sup>, MassBank of Japan (<http://massbank.jp>)<sup>74</sup>, EU  
305 (<https://massbank.eu/MassBank/>)<sup>75</sup> and North America  
306 (<http://mona.fiehnlab.ucdavis.edu/>), mzCloud<sup>76, 77</sup>, and ReSpec<sup>78</sup>, which potentially  
307 provides users with access to around 2.4 million MS<sup>2</sup> reference spectra, when GC-MS and  
308 LC-MS reference spectra are both considered<sup>66</sup>. Many of these reference databases have  
309 an integrated spectral matching tool for compound identification, including mzCloud,  
310 METLIN/XCMS Online<sup>79, 80</sup>, Metabox<sup>81</sup>, MassBank). The goal of GNPS is not only to  
311 provide a spectral matching tool, but also to serve as a data storage and knowledge  
312 capture and dissemination platform, and to provide access to a host of other analysis tools  
313 not covered in detail here, such as *in silico*-based dereplication<sup>82-84</sup>, network annotation  
314 propagation<sup>85</sup>, genome mining tools<sup>86</sup>, and MASST searches.

315 GNPS is currently the only online platform that provides molecular networking, a  
316 computational tool that compares pairs of MS<sup>2</sup> spectra based on their similarities and  
317 connects them to MS<sup>2</sup> reference spectral libraries. Molecular networking enables further  
318 propagation of annotations through mass spectral relations. MetGem<sup>87</sup> is a standalone  
319 software package that can be used for the generation of molecular networks which works  
320 well for smaller data sets, it is not connected to a knowledge base, repository wide analysis  
321 tools and additional computational resources that GNPS provides.

322

#### 323 **1.4 Expertise needed to implement the protocol**

324 Sampling and sample preparation, including sample extraction, should be  
325 performed by a trained analytical chemist, and mass spectrometry data should be acquired  
326 by a trained mass spectrometrist. It is imperative that the parameters for mass  
327 spectrometry be suitably optimized for the experimental conditions and sample type in  
328 order to generate meaningful molecular networks. Important instrument parameters to  
329 consider may include precursor isolation window, mass resolution, collision energy, data  
330 dependent acquisition settings (e.g. duty cycle time and dynamic exclusion parameters),  
331 and the mass spectrometer has to be properly calibrated before use. While an expert user  
332 will have preferred instrument parameters, recommended data acquisition parameters  
333 from major instrument manufacturers are provided below (section 2.6) for newer mass  
334 spectrometry users who aim to create molecular networks in GNPS. Basic knowledge of  
335 tandem mass spectrometry fundamentals as well as knowledge of sample handling and  
336 preparation are required to further optimize the data analysis parameters appropriate to  
337 the instrument used and the experimental design.

338

#### 339 **1.5 Experimental design**

340 After running the molecular networking algorithm, GNPS creates a data table that  
341 provides as much chemical insight into the data as possible in relation to the metadata  
342 (associated sample information) provided by the user. Such data tables can be viewed as  
343 networks directly in the GNPS website or exported and manipulated in other data  
344 visualization tools and statistical analysis packages. Here we provide a GNPS-based  
345 molecular networking tutorial in which we import the table into a third party tool called  
346 Cytoscape, a powerful network visualization software. Notably, the information  
347 represented in and inferred from a molecular network is dependent on the input, including  
348 both the mass spectrometry data<sup>88</sup> and networking parameters selected.

349

##### 350 **1.5.1 Reproducibility, blanks, and controls**

351 A well organized and well thought-out experimental plan is essential for the  
352 successful creation of useful molecular networks, since molecular networks are only as  
353 meaningful as the experiment and data from which they originate. This includes providing  
354 sample information (metadata) tables and raw data files for the sample set; metadata  
355 tables aid the creation of molecular networks that have increased interpretative value. In  
356 order to avoid pitfalls associated with large-scale mass spectrometry experiments, e.g.  
357 batch effects<sup>89</sup>, sample carryover and/or contamination<sup>90</sup>, and high background signal<sup>91</sup>,  
358 and to maximize reproducibility and signal-to-noise ratio<sup>92</sup>, a dataset should include  
359 blanks, quality control (QC) samples, and experimental replicates. Dunn et al.<sup>93</sup> describe  
360 an appropriate representative experimental design in detail that includes blanks, quality  
361 control mixtures, and samples plus internal standards. Petras et. al.<sup>38</sup> provide an example  
362 that illustrates control metrics, including evaluation of quality control mixtures and signal  
363 deviation of the internal standard.

364 We recommend preparing control samples using exactly the same protocols and  
365 experimental conditions used to prepare test samples (i.e. the same types of tubes, the  
366 same batches of tubes, the same extraction solvent, extraction time, sonication time/power  
367 and so on). These blank samples inform which ions come from the experimental conditions  
368 and they can be subtracted from test sample signals in the molecular networking analysis  
369 (see section 3.4.3). The requirements for QC associated with a broad assessment of the  
370 natural product composition of an extract library used in bioactivity screens is different from  
371 a detailed clinical study for biomarker discovery. When possible, one should add internal  
372 standard(s) to each sample to ensure that the system performs consistently. If the internal  
373 standard(s) do not match the user-defined acceptable chromatography variations, the  
374 sample needs to be either removed from downstream analysis or rerun. This is particularly  
375 useful in applications where thousands of samples, such as natural product extract  
376 libraries, are screened. Further, when acquiring data for a large number of samples,  
377 especially when multiple batches are used, we suggest acquiring data for additional QC  
378 samples to monitor batch and plate effects throughout the experiment in order to assess  
379 instrumental variations over time, such as retention time drift. QC samples may either  
380 consist of aliquots from a subset of test samples pooled together (pooled QC) or be  
381 mixtures of molecules specifically defined for quality assurance. For example, it is common  
382 to use the last column of a 96-well plate for the QC mixture to ensure that the instrument  
383 and chromatography behave in an identical fashion throughout an experiment. Finally,  
384 data from experimental replicates, including both technical and biological replicates should  
385 be acquired in a randomized fashion. This is especially important for large-scale population  
386 studies to ensure minimized bias. One common problem in metabolomics and LC-MS  
387 analysis is sample carryover, caused by residual compound(s) from a previous run. One  
388 way to reduce this issue is to insert a wash routine between samples followed by a blank  
389 to ensure that no carryover is observed.

### 391 **1.5.2 Molecular networking parameters**

392 GNPS-based molecular networking parameters may be varied significantly and need to  
393 be set appropriately for the acquired dataset, based on sample (anticipated molecular  
394 masses and types of molecules), instrument resolution and collision energies used for MS  
395 acquisition. Networking parameters are described in detail in Section 3.3, Table 1, and  
396 should be considered and selected carefully in order to obtain useful networks, which  
397 ultimately depend on the quality and quantity of MS<sup>2</sup> spectra.

398

## 399 1.6 Limitations and challenges

400 Since GNPS-based molecular networking utilizes MS<sup>2</sup> data, it is susceptible to the  
401 same challenges encountered in any mass spectrometry data acquisition experiment,  
402 such as low signal-to-noise, insufficient separation of analytes, or poor peak shape.<sup>94, 95</sup> In  
403 addition, classical molecular networking can provide only qualitative information about the  
404 experiment because only MS<sup>2</sup> scans are considered in the analysis. While feature-based  
405 molecular networking (Box 4) incorporates MS<sup>1</sup> and chromatographic data, which  
406 approximates quantitation, it is still not strictly quantitative. If calibrated quantitative  
407 information is needed to answer the scientific question, follow-up experiments should be  
408 performed using targeted LC-MS.

409 Additionally, one should consider potential issues that accompany metabolomics  
410 experiments, such as sample extraction efficiency and reproducibility, as well as unwanted  
411 metabolite degradation. While avoiding degradation or modification of all molecules in a  
412 sample is impossible, it is important that all samples for comparison are prepared and  
413 analyzed in an identical manner, unless the goal is to understand the effects of sample  
414 preparation conditions<sup>96</sup>. While a few publications describe the impact of storage on the  
415 detectable metabolome, these are sample type-specific and there is currently no  
416 consensus for a “gold standard”<sup>97-99</sup>. Ultimately, sample preparation is highly dependent  
417 on the type of sample collected, and includes drying, homogenization, and extraction  
418 steps<sup>100</sup>. Although every lab has their own preferences for sample treatment, we strongly  
419 advocate for samples to be collected and extracted with solvent as soon as possible. The  
420 speed of this is dependent on the experimental environment. For example, samples  
421 collected in remote areas, at sea using a small boats, or often even in a clinical setting,  
422 may be stored for hours or days before they can be extracted, given that some solvents  
423 are not easily brought into a clinical setting or used while out at sea. In contrast, samples  
424 from a cultured system in a lab or an enzymatic reaction, for example, can be halted in  
425 milliseconds using a rapid quench system and can then be extracted in seconds. The  
426 choice of solvent and extraction protocol is dictated by the experimentalist’s interests and  
427 questions. Although there is always overlap among the molecules from even very different  
428 extraction protocols, more polar metabolites are extracted with ethanol, methanol and  
429 butanol while more hydrophobic metabolites are extracted with benzene, ethyl acetate or  
430 chloroform<sup>96</sup>. The samples can then be introduced into the mass spectrometer using front-  
431 end separation techniques, most often liquid chromatography or ion mobility. If mass  
432 spectrometry cannot be performed immediately, we recommend completely drying the  
433 samples before storage at cryogenic temperatures.

434 To annotate unknown molecules, GNPS queries MS<sup>2</sup> spectra against MS<sup>2</sup> data in  
435 reference libraries and assigns a cosine score based on their similarity. For the GNPS  
436 spectral library, MS<sup>2</sup> spectra are acquired from laboratories around the world using a  
437 variety of mass spectrometers and sample preparation protocols. Therefore, mass spectra  
438 submitted to GNPS can differ in terms of both quality and content. For instance, MS<sup>2</sup>  
439 fragment ions and their intensities can vary significantly between instruments, and even  
440 on the same instrument if the experimental setup is changed<sup>101</sup>. GNPS requires that the  
441 instrument and ion source be specified with each reference spectrum submitted and it is  
442 recommended that this be taken into account when assessing the quality of a library hit.  
443 Along these lines, annotations of unknown molecules are not all accurate and should be  
444 considered putative until confirmed with an authentic chemical standard.

445 On average, in 2016 when GNPS was published, only 2% of spectra in an  
446 untargeted mass spectrometry metabolomics experiment were annotated<sup>102</sup>. Although this

447 percentage has grown to an average of 5-6% annotations, a large percentage of MS<sup>2</sup>  
448 spectra typically remain unannotated. The structures of these unannotated molecules or  
449 “dark matter”<sup>103</sup> might be known, but their identity is not revealed because no reference  
450 spectra exist in library databases, against which to compare. To improve annotation rates,  
451 *in silico* tools have been developed to match unknown MS<sup>2</sup> spectra to putative chemical  
452 structures<sup>104</sup>. Several of these computational tools, which include MetFrag<sup>105</sup>,  
453 MetFusion<sup>106</sup>, SIRIUS<sup>107, 108</sup>, CSI:FingerID<sup>109</sup>, MS-Finder<sup>110</sup>, Network Annotation  
454 Propagation (NAP)<sup>85</sup>, and Dereplicator<sup>82, 83</sup> can be integrated into GNPS molecular  
455 networking workflows to provide insight into the annotation; the application of such tools is  
456 beyond the immediate scope of the networking protocol presented here.

457

## 458 2.0 Materials

### 459 2.1 REAGENTS

460 **CRITICAL** For specific storage and handling instructions, consult the manufacturer of each  
461 reagent. Although high grade solvents are used, different batches of the same solvents  
462 (even purchased from the same vendors), can give rise to different background  
463 contaminants in the experiment. There are also many possible substitutes for the reagents  
464 and consumables listed below.

- 465 • Water of LC–MS (Optima) grade (Thermo Fisher Scientific, cat. no. W6-4)
- 466 • Acetonitrile (ACN), LC–MS (Optima) grade (Thermo Fisher Scientific, catalogue  
467 number A955-4) ! **CAUTION** Acetonitrile is highly flammable, and the bottles  
468 should be stored in a flammable-liquid cabinet.
- 469 • Methanol (MeOH), LC-MS grade ! **CAUTION** Methanol is highly flammable, and  
470 the bottles should be stored in a flammable-liquid cabinet.
- 471 • Formic acid (FA). LC–MS grade, Optima grade (Thermo Fisher Scientific,  
472 catalogue number A117-50) ! **CAUTION** Formic acid is highly corrosive. It should  
473 be handled in a flow cabinet while wearing eye protection and gloves.
- 474 • LC-MS calibration solutions, e.g. for the Bruker MaXis II QTOF mass spectrometer:  
475 ESI-TOF Low Concentration Tuning Mix (Agilent Technologies, catalogue number  
476 G1969-85000) for external calibration and Hexakis(1H,1H,3H-  
477 tetrafluoropropoxy)phosphazene (Synquest Laboratories, catalogue number  
478 8H79-3-08), *m/z* 922.009798 for internal calibration (lock mass) ! **CAUTION** This  
479 compound is irritating to the eyes and the skin. It should be handled wearing eye  
480 protection and gloves; for the Q-Exactive mass spectrometer: Pierce LTQ Velos  
481 ESI Positive Ion Calibration Solution (Thermo Fisher Scientific, catalogue number  
482 88323) and ESI Negative Ion Calibration Solution (Thermo Fisher Scientific,  
483 catalogue number 88324).

484

### 485 2.2 EQUIPMENT

- 486 • Microtiter plates (e.g. Nunc 96-Well Round Bottom Polypropylene Storage  
487 Microplates, Thermo Fisher Scientific, catalogue number 267245) containing  
488 samples of interest at, e.g., 1 mg/mL concentration.
- 489 • Benchtop vacuum concentrator compatible with 96-well microplate evaporation  
490 (Centrivap; Labconco)
- 491 • Reversed phase C18 LC column, 1.7- $\mu$ m particle size, 50  $\times$  2.1-mm (Phenomenex,  
492 part number 00B-4475-AN or equivalent)
- 493 • UHPLC system coupled to a tandem mass spectrometer with an ESI source; e.g.  
494 a 1260 HPLC (Agilent) coupled to a QTOF 6530 mass spectrometer (Agilent),

495 UltiMate 3000 UHPLC system (Dionex) coupled to a MaXis II QTOF system  
496 (Bruker Daltonics), a Vanquish UHPLC system coupled to a Q-Exactive mass  
497 spectrometer (Thermo Fisher Scientific), an Acquity UHPLC I coupled to a Xevo  
498 G2-XS QTOF (Waters), a Nexera X2 UHPLC (or a Prominence UFLC) coupled to  
499 an IT-TOF mass spectrometer (Shimadzu) or an AB Sciex 5600 TripleTOF mass  
500 spectrometer.

501

## 502 2.3 SOFTWARE

- 503 • MSConvert tool from the ProteoWizard  
504 (<http://proteowizard.sourceforge.net/downloads.shtml>)
- 505 • AB Sciex MS Data Converter (Beta 1.3) is freely available for download from the  
506 AB Sciex website <https://sciex.com/software-support/software-downloads>
- 507 • AB Sciex Analyst Software 1.7 is available for download, trial license use and  
508 purchase from the AB Sciex website [https://sciex.com/products/software/analyst-](https://sciex.com/products/software/analyst-software)  
509 [software](https://sciex.com/products/software/analyst-software)
- 510 • Agilent MassHunter software can be obtained from the Agilent website:  
511 [https://www.agilent.com/en/products/software-informatics/masshunter-](https://www.agilent.com/en/products/software-informatics/masshunter-suite/masshunter/masshunter-software)  
512 [suite/masshunter/masshunter-software](https://www.agilent.com/en/products/software-informatics/masshunter-suite/masshunter/masshunter-software)
- 513 • Bruker DataAnalysis is available for download from the Bruker website  
514 ([www.bruker.com/service/support-upgrades/software-downloads/mass-](http://www.bruker.com/service/support-upgrades/software-downloads/mass-spectrometry.html)  
515 [spectrometry.html](http://www.bruker.com/service/support-upgrades/software-downloads/mass-spectrometry.html))
- 516 • Shimadzu LabSolutions can be obtained from the Shimadzu website:  
517 [https://www.ssi.shimadzu.com/products/liquid-chromatography-mass-](https://www.ssi.shimadzu.com/products/liquid-chromatography-mass-spectrometry/lcms-software.html)  
518 [spectrometry/lcms-software.html](https://www.ssi.shimadzu.com/products/liquid-chromatography-mass-spectrometry/lcms-software.html)
- 519 • Thermo Scientific Xcalibur software can be obtained at:  
520 <https://www.thermofisher.com/order/catalog/product/OPTON-30801>
- 521 • Waters MassLynx MS software can be obtained at:  
522 [http://www.waters.com/waters/en\\_US/MassLynx-MS-](http://www.waters.com/waters/en_US/MassLynx-MS-Software/nav.htm?locale=en_US&cid=513662)  
523 [Software/nav.htm?locale=en\\_US&cid=513662](http://www.waters.com/waters/en_US/MassLynx-MS-Software/nav.htm?locale=en_US&cid=513662)
- 524 • FTP Client (e.g. WinSCP for Windows; Cyberduck for Macintosh)
- 525 • Web Browser, Firefox or Google Chrome to access GNPS
- 526 • Cytoscape for data visualization: <https://cytoscape.org/> ([current version at the time](https://cytoscape.org/)  
527 [of publication is 3.7.1](https://cytoscape.org/)).
- 528 • Software relevant to optional pipelines, e.g. 2D or 3D Visualization<sup>111</sup>; Feature-  
529 based molecular networking, see **Box 4**.

530

## 531 2.4 EXAMPLE DATASETS

532 **CRITICAL** All LC–MS data used in this paper are publicly available at the GNPS-MassIVE  
533 repository under the following accession numbers.

- 534 • [MSV000083437](https://massive.ucsf.edu/ncgi/study/MSV000083437) (Germ Free and Specific Pathogen Free Mice, unpublished)
- 535 • [MSV000083359](https://massive.ucsf.edu/ncgi/study/MSV000083359) (3D Cartography of Diseased Human Lung<sup>50</sup>)
- 536 • [MSV000083381](https://massive.ucsf.edu/ncgi/study/MSV000083381) (Stenothricin-GNPS analogues<sup>11</sup>)

537

## 538 2.5 REAGENT SETUP

539 **Aqueous LC–MS mobile phase, Solvent A** Prepare the aqueous mobile phase (Solvent  
540 A) for LC–MS by adding LC–MS-grade formic acid to LC–MS-grade water to make a  
541 100:0.1 (vol/vol) water/formic acid mixture. The LC solvents can be stored at room  
542 temperature for up to 1 week.

543 **! CAUTION** Formic acid is highly corrosive. **CRITICAL** The aqueous mobile phase for LC–  
544 MS should not be stored for more than a week because of the potential for microbial  
545 growth.

546 **Organic LC–MS mobile phase, Solvent B** Prepare the organic mobile phase (Solvent B)  
547 for LC–MS by adding LC–MS-grade formic acid to LC–MS-grade acetonitrile to make a  
548 100:0.1 (vol/vol) acetonitrile/formic acid mixture.

549 The LC solvents can be stored at room temperature for up to 1 week.

550 **! CAUTION** Formic acid is highly corrosive and should be handled in a flow cabinet while  
551 wearing eye protection and gloves.

552

## 553 **2.6 EQUIPMENT SETUP**

### 554 **Mass spectrometry**

555 Both ion source parameters and data dependent acquisition (DDA) parameters are  
556 essential for obtaining quality MS<sup>2</sup> spectra to be used for molecular networking. Although  
557 many instrument configurations exist, several representative ion source and DDA  
558 parameters are described below. Relevant to these MS parameters is the LC method used,  
559 an example of which is a gradient profile from 10 to 100% ACN + 0.1% FA in H<sub>2</sub>O + 0.1%  
560 FA (for 12 min), followed by isocratic 100% ACN + 0.1% FA (for 3 min), and 5% ACN +  
561 0.1% FA (3 min) re-equilibration phase, with a flow rate of 400 µL/min.

562

563 Suggested instrument parameters for ABSciex, Agilent, Bruker, Shimadzu, Thermo  
564 Scientific, and Waters are provided in the supporting information.

565

### 566 **3.0 Procedure**

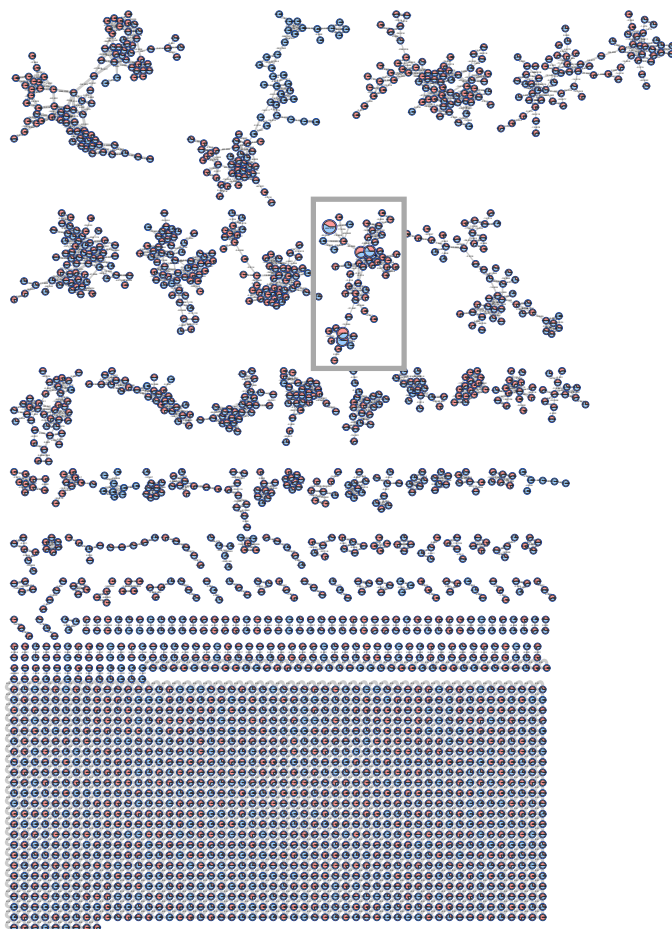
567

568 In addition to the protocol described in the following, all steps, albeit in less detail, are also  
569 described and continuously updated and maintained in the online GNPS documentation  
570 at: <https://ccms-ucsd.github.io/GNPSDocumentation/>

571

572 The data submission and molecular networking workflow (section 3.2 onwards) may be  
573 followed as a tutorial using an untargeted metabolomics dataset for 3D molecular  
574 cartography of the mouse duodenum (paper is in review, Massive dataset  
575 [MSV000083437](https://doi.org/10.26434/chemrxiv-2023-08343)). This dataset is a subset of a collection of metabolomes analysed from  
576 organs of germ free (GF) and specific pathogen free (SPF) mice that led to the discovery  
577 of new amide conjugated bile acids made by bacteria that affect host metabolism *via*  
578 farnesoid X receptor (FXR) agonism. The following procedure will take the reader through  
579 submission of dataset [MSV000083437](https://doi.org/10.26434/chemrxiv-2023-08343) to the molecular networking workflow in GNPS,  
580 through the molecular networking workflow in GNPS (including input parameters), and  
581 through visualization of the generated network using both in browser and Cytoscape-  
582 based visualization (**Fig. 3**).





- = Germ free (GF) mice
- = Specific Pathogen Free (SPF) mice

583  
584

585 **Figure 3.** The readers will recreate the mouse duodenum global molecular network  
586 depicted above, created from MassIVE dataset [MSV000083437](https://massive.ucsd.edu/MSV000083437) and visualized in  
587 Cytoscape. Pie charts represent relative summed precursor ion intensities per MS<sup>2</sup> spectra  
588 detected within each metadata group: red for germ-free (GF) and blue for specific-  
589 pathogen free (SPF) mice. The Box highlights a cluster we will examine below in terms of  
590 chemical interpretation.

591

592 **3.1 Data conversion** - Timing 1 hour up to a few days (varies depending on size of dataset  
593 and computer set-up)

594

595 GNPS-MassIVE converts raw data formats after upload to .mzML format (stored in the  
596 ccms\_peak folder) for GNPS processing. Nevertheless, to enable immediate use of the  
597 data, it is recommended to manually convert the raw data to open file formats prior to  
598 uploading to GNPS-MassIVE. The protocol for data conversion depends on the instrument  
599 used for mass spectrometry acquisition. MSConvert can be used for the conversion to a  
600 GNPS-compatible format of mass spectrometry data acquired on AB Sciex, Agilent,  
601 Shimadzu (after initial conversion, SI Methods), Thermo Scientific and Waters instruments.  
602 Although one of the most common formats used in GNPS, Bruker files (.d format), at this  
603 time, are still not MSConvert compatible. For Bruker files, a separate workflow must be

604 utilized, which applies internal lockmass calibration to the output file. This Bruker workflow  
605 is described in more detail in the SI methods. Alternatively, for AB Sciex, raw files (.wiff)  
606 could be converted into .mzML format using the AB MS Data Converter (AB Sciex version  
607 1.3 beta, freely available at <https://sciex.com/software-support/software-downloads>).

608

609 1) MSConvert can be downloaded freely from ProteoWizard at:  
610 <http://proteowizard.sourceforge.net/download.html>. This software is compatible with  
611 Windows and Linux operating systems but is not supported for Mac OS. When  
612 downloading ProteoWizard, the version of Windows must be specified and .NET  
613 Framework 3.5 SP1 and 4.0 must be installed. Then either a traditional workflow or an  
614 easy workflow can be used for the file conversion. These two workflows are detailed below.  
615 The “traditional” workflow, outlined below, is the manual workflow.

616

617 2) Mass spectrometry files must be converted to open file formats such as .mzXML,  
618 .mzML, and .mgf formats for analysis in GNPS, with the preferred formats being .mzXML  
619 and .mzML. Although it is encouraged to co-submit the raw data to MassIVE, GNPS does  
620 not support .mzData, .xml, .raw, .wiff, .scan, .d, and .cdf formats.

621

622 3) MSConvert is the recommended software for conversion of data acquired on AB Sciex,  
623 Agilent, Thermo Scientific and Waters instruments. Conversion can be performed following  
624 the steps outlined below:

625 a. In the Start Menu, the ProteoWizard folder can be selected and MSConvert can be  
626 opened.

627 b. To select file(s) for conversion, click Browse; then click ‘Add’ to add file(s) to the  
628 workflow and select a directory for the output.

629 c. To convert the vendor file format to an .mzXML file, select .mzXML under Options;  
630 32-bit should be selected for binary encoding precision and Use zlib compression  
631 should be unchecked.

632 d. Choose Peak Picking under the Filters heading and under Algorithm check Vendor,  
633 then write in MS-Levels 1-2 and finally add the filter by clicking Add. **! CRITICAL**  
634 **STEP** Move the peakPicking filter to the top of filter list. The peakPicking filter must  
635 be the first filter in the list or the output file will not be centroided.

636 e. Click Start then check the folder for the .mzXML files in the Output Directory. These  
637 files can be opened in SeeMS (Installed with MSConvert), OpenMS TOPPView  
638 (<https://github.com/OpenMS/OpenMS/releases>)<sup>112</sup> or MZmine2  
639 (<https://github.com/mzmine/mzmine2/releases>)<sup>113</sup> to verify that the conversion  
640 worked properly.

641

642 An “easy” workflow is also available. This simple batch conversion method includes a  
643 complete package for Windows users to convert vendor formats to GNPS compatible  
644 format (mzXML, mzML, MGF) and is described in the SI Methods. An online data  
645 conversion tutorial can be accessed at: [https://ccms-](https://ccms-ucsd.github.io/GNPSDocumentation/fileconversion/)  
646 [ucsd.github.io/GNPSDocumentation/fileconversion/](https://ccms-ucsd.github.io/GNPSDocumentation/fileconversion/).

647

### 648 **3.2 Data submission to GNPS / MassIVE**

649 It is necessary to create an account with GNPS in order to submit datasets and create  
650 workflows, as well as to receive emails about the outcomes. Making a GNPS account  
651 automatically sets up a MassIVE account that uses the same login and password. To

652 manipulate MS data files in GNPS, they must first be uploaded to MassIVE, which is an  
653 online repository for mass spectrometry datasets hosted by the UCSD Center for  
654 Computational Mass Spectrometry (CCMS). The user workspace in GNPS / MassIVE  
655 provides a personalized location for researchers to curate mass datasets, submit and  
656 monitor GNPS workflows, subscriptions to datasets that have been made publicly  
657 available by others, or clone and reanalyze either their own or other public datasets. More  
658 information on subscriptions to data can be found in sections 3.6.

659

### 660 **3.2.1 Create a GNPS / MassIVE account (SI Fig. 1):**

661

662 1) Open up a web browser. GNPS is designed to work with Firefox or Google Chrome  
663 but also works in Microsoft Edge, Safari, and Opera.

664 2) Navigate to the GNPS home page by using this link  
665 <https://gnps.ucsd.edu/ProteoSAFe/static/gnps-splash2.jsp>

666 3) Towards the top center of the page, above the large GNPS logo, click on “Register  
667 New Account” (right hand grey box).

668 4) On the new page that loads, enter a username, name (optional), organization  
669 (optional), email, and password (twice for confirmation) in the spaces provided.

670 4) Click submit.

671 5) Sign-in to your new GNPS account <http://massive.ucsd.edu/ProteoSAFe/> and  
672 check that your GNPS credentials work for logging in to MassIVE.

673

### 674 **Box 1: Navigating the User Workspace Portal**



675

676 At the top of the GNPS website, users will find a banner that allows them to navigate their  
677 personal workspace and access additional resources such as the help forum and  
678 molecular networking documentation. Within this space, the ‘My User’ tab provides a way  
679 to view all MassIVE datasets and reference spectra deposited by the user, and the ‘Jobs’  
680 button allows easy access to all jobs submitted by the user through the GNPS and  
681 MassIVE interfaces. Clicking on ‘MassIVE datasets’ allows the user to browse and  
682 subscribe (section 3.7.3) to all public MassIVE datasets with GNPS in the title. Additionally,  
683 this banner is a portal to all resources for help using GNPS. The ‘Documentation’ link in  
684 the banner takes the user to the GNPS documentation website, which has step-by-step  
685 instructions and links to tutorial videos as well as access to the ‘legacy’ documentation  
686 (from a menu on the right-hand side of the page) that can provide additional information  
687 to the user. The ‘Forum’ button opens a Google groups forum where users can post  
688 questions, have discussions and report potential bugs. The corresponding online tutorial  
689 can be accessed at: <https://ccms-ucsd.github.io/GNPSDocumentation/quickstart/>

690

### 691 **3.2.2 Deposit data files by submitting a dataset (SI Fig. 3) :**

692 An online tutorial on how to submit a dataset to MassIVE can be accessed at: <https://ccms-ucsd.github.io/GNPSDocumentation/datasets/#submitting-gnps-massive-datasets>

693

694

695

There are two steps to submitting a dataset to the GNPS-MassIVE repository:

696 **Step 1 (SI Fig. 3a).** Upload your data files to the MassIVE web server using an  
697 FTP client - Timing 10 min to get the upload process started.

698 Of the many free dedicated FTP clients, the following are more popular ones that have  
699 been tested with MassIVE: WinSCP, CoreFTP, and CoffeeCup Free FTP for Windows,  
700 and Cyberduck or FileZilla for Macintosh. Caution: when downloading an FTP client for  
701 use, make sure it comes from a trusted source to avoid malware. Data files transferred to  
702 MassIVE should be in .mzXML, .mzML, .mgf formats. The data that is uploaded should  
703 **not be in a file archive (e.g. zip, tar) format.** It is also encouraged that the original vendor  
704 raw data files (e.g. .wiff for AB Sciex, .yep for Agilent, .d for Bruker, .lcd for Shimadzu, .raw  
705 for Thermo Scientific) are uploaded together with the open formats as described below.

706  
707 1) Change file protocol to FTP and log onto the FTP server with the host name  
708 *massive.ucsd.edu* using your MassIVE web account username and password in the FTP  
709 client program for FTP file transfer. Most FTP clients use this "Quick Connect" feature.  
710 Alternatively, type in the FTP server name, username and password, and then connect  
711 directly.

712  
713 **Step 2 (SI Fig. 3b).** Run the MassIVE dataset submission workflow on the  
714 uploaded files as follows:

715  
716 1) Load the home page for MassIVE from the GNPS home page by scrolling down to  
717 the GNPS-MassIVE datasets section and click on the 'Deposit dataset' bar in the 'Create  
718 Public datasets' block. Alternatively, click on the 'Submit your data' link in the paragraph  
719 titled 'Submit Data' on the MassIVE home page. A direct way to deposit the data is to  
720 navigate directly to the MassIVE home page (<http://massive.ucsd.edu/ProteoSAFe/>). This  
721 will bring up the Dataset Submission workflow input form, on which there are varying  
722 numbers of fillable fields under each of the following sections described below.

723  
724 The reader can follow along (**SI Fig. 3**), as this has already been completed for the  
725 MassIVE dataset [MSV000083437](#).

726  
727 2) In the 'Workflow Selection' section:  
728 Enter a title for your dataset, **noting that GNPS datasets must have a 'GNPS'**  
729 **prefix in the title** in order for these GNPS-MassIVE datasets to be visible to GNPS users.  
730 **Adding GNPS in the title is therefore absolutely !IMPORTANT! for the dataset to**  
731 **become a part of the community and ensures that the data becomes alive (Section**  
732 **3.6) and enables subscriptions and other analysis features specifically used for the**  
733 **GNPS community (Section 3.6).** If a "GNPS" tag is not added at the beginning of the title  
734 it will not be part of the GNPS analysis infrastructure. Currently all of MassIVE has almost  
735 ~11,000 public mass spectrometry datasets (mostly proteomics), ~1,100 of which are also  
736 part of GNPS. If GNPS is not added from the beginning it is possible to go to MassIVE,  
737 log-in and edit the title at a later time.

738  
739 To satisfy this requirement for the dataset that reader will use in this tutorial,  
740 MassIVE dataset [MSV000083437](#) has been titled "GNPS Example Dataset\_GF vs. SPF  
741 Mouse Duodenum."  
742

743 3) In the 'Dataset Metadata' section:

744 To minimize the burden to make datasets for GNPS analysis and to enable as  
 745 much flexibility in what additional information the user wants to make available, very few  
 746 metadata fields are absolutely required, although the user is encouraged to provide as  
 747 much metadata as possible. It should be noted that the datasets that have the most  
 748 information associated with it are also the datasets that are the most visible to the  
 749 community. Fields for metadata relevant to the dataset being submitted are listed in the  
 750 table below. The first three fields ('Species', 'Instrument' and 'Post-Translational  
 751 Modifications') are backed by lists of standardized controlled vocabulary (CV) terms,  
 752 maintained by organizations such as the [HUPO Proteomics Standards Initiative](#)<sup>114</sup> and  
 753 many others CVs that the user can implement<sup>114, 115</sup>. To search these terms, type at least  
 754 3 characters into any of these text boxes, and a drop-down list of supported terms that  
 755 match your query will be displayed. To select a term, click on it in the drop-down list and it  
 756 will be added to your dataset. **Using the official CV to tag your dataset greatly  
 757 increases the likelihood that it will be found and processed correctly by any  
 758 automated software that may interface with the MassIVE repository.** If the term you  
 759 want is not present in the list, you can type your custom text in the text box and click the  
 760 adjacent 'Add' button to tag your dataset.

761

762 **Table 2.** Metadata Categories for Data Upload to MassIVE

Metadata Category	Required	Notes	Example Dataset <a href="#">MSV000083437</a>
Species	Yes	Enter custom text if the correct species for your dataset is not supported in the list or if you sample is not a specific species (e.g. environmental sample or community of organisms).	<i>Mus musculus</i> (house mouse)
Instrument	Yes	Enter custom text if the correct instrument for your dataset is not supported in the list.	maXis
Post-Translational Modifications	Yes	For small molecule metabolomics datasets the appropriate entry in the drop-down list is: 'PRIDE:0000398, No PTMs are included in the dataset'.	No PTMs included in the dataset
Keywords to assign to your dataset	Yes	Your dataset must be tagged with at least one keyword - there is no limit. Keywords are custom text, so you must click the 'Add' button after entering text.	mouse duodenum
Principal	Yes	To identify the lab providing	Pieter Dorrestein

Investigator		the data.	( <a href="mailto:pdorrestein@ucsd.edu">pdorrestein@ucsd.edu</a> ) UCSD, United States
Description	No	Recommended to provide as much detail as possible	N/A

763

764 Metadata (sample information) for MassIVE dataset [MSV000083437](#) has been added as  
765 shown in **SI Fig. 2b** and is tabulated above.

766

767 4) In the 'Dataset File Selection' section there are eleven different file types that can  
768 be added and these are organized into three different categories - required, recommended  
769 or optional. **Most of these file categories are not strictly required. The only official**  
770 **file requirement for a MassIVE dataset is that at least one file is submitted in either**  
771 **the 'Raw Spectrum Files' or 'Peak List Files' categories. If a submitted dataset does**  
772 **not meet the additional requirements for a '[complete](#)' submission, then it is**  
773 **considered 'partial', which is currently standard for small molecule datasets that are**  
774 **a part of GNPS.**

775

a) *Recommended for all submissions*

776

i) Raw Spectrum Files – Raw mass spectrum files in a non-standard or  
777 instrument-specific format, such as AB Sciex .wiff files, Agilent .yep files,  
778 Shimadzu .lcd files, Bruker .d files Thermo Scientific .raw files, Waters .raw  
779 files.

780

ii) Peak List Files – Processed mass spectrum files in a standardized format.  
781 The following formats are recognized by MassIVE as valid for this category:  
782 .mzXML, .mzML, and .mgf. This is the file from which GNPS analysis is  
783 enabled.

784

b) *Strongly encouraged for submissions to improve the ability to interpret the final  
785 molecular networks.*

786

i) Supplementary Files – All remaining files relevant to this dataset that do not  
787 properly fit into any of the other listed file categories. **A metadata file**  
788 **(sample information in a tab delimited text format) with relevant**  
789 **attributes that can be used for visualizing the data in networks should**  
790 **be included here (see Box 3).**

791

c) *Required for "Complete" Submission* Result Files – **Not necessary for small  
792 molecule workflows - although possible and encouraged.** Spectrum  
793 identifications in a standardized format. The following formats are recognized by  
794 MassIVE as valid for this category: mzIdentML<sup>116</sup> and mzTab<sup>117</sup>, mzTab-M<sup>118</sup>.

795

i) Search Engine Files – The output of any search engine or data analysis  
796 tools or pipelines that were used to analyze this dataset, unless provided in  
797 a standardized format recognized by the 'Result Files' category (see  
798 above).

799

d) *Optional*

800

i) License Files – Specifying how and under what conditions the dataset files  
801 may be downloaded and used. Multiple license files may be uploaded, if  
802 appropriate. By default, you can simply leave the 'Standard License'  
803 checkbox checked and your dataset will be submitted under the default  
804 [Creative Commons CC0 1.0 Universal](#) license. However, if you wish to

- 805 provide your own license, then you can uncheck this box and then assign  
806 your own file to the 'License Files' category.
- 807 ii) Spectral Libraries – Any custom spectral library files that were searched  
808 against in the analysis of this dataset, or that were generated using the  
809 spectrum files provided in this dataset, if applicable.
  - 810 iii) Methods and Protocols – Any open-format files containing explanations or  
811 discussions of the experimental procedures used to obtain or analyze this  
812 dataset.
- 813 e) *Optional, mostly relevant to peptidomics and proteomics projects*
- 814 i) Quantification Results – Any data and metadata generated by the analysis  
815 software used. Typically applied to the quantification analysis of peptides  
816 and proteins.
  - 817 ii) Gel Images – Any gel image files generated, in the event that two-  
818 dimensional gel electrophoresis has been used as a separation method.
  - 819 iii) Sequence Databases – Any files from protein or other sequence databases  
820 that were associated with or searched against in the analysis of this dataset,  
821 if applicable (usually .fasta format).
- 822

823 For readers that are following the example, peak List files were uploaded previously for  
824 dataset [MSV000083437](#), as illustrated in **SI Fig 3b**, where nine folders (Control, GF1,  
825 GF2, GF3, GF4, SPF1, SPF2, SPF3, SPF4) have been added.

826

827 5) 'Mapping Spectrum Files to Identification Files' is **not necessary for small**  
828 **molecule workflows**. In order for a submission to qualify as 'complete', each spectrum  
829 (data) file referenced within a "Result File" must be associated with a file from the "Peak  
830 List Files" category. This section is where these two types of files are associated with each  
831 other as appropriate.

832

833 6) The 'Dataset Publication' section has three optional fields to:

- 834 a) 'Enter a Password' (e.g. to share selectively with collaborators and  
835 manuscript reviewers),
- 836 b) 'Share on ProteomeXchange' is **not applicable to small molecule**  
837 **workflows**: checking the box will submit and announce the dataset via the  
838 ProteomeXchange consortium at the time that it is made public on  
839 MassIVE. The dataset will not appear publicly in either repository until you  
840 click the 'Make Public' button on your dataset's status page (see below).
- 841 c) 'Generate a DOI' if you want a Digital Object Identifier to be generated and  
842 assigned to this dataset. This is encouraged for all public datasets and can  
843 be used in publications.

844

845 7) The section titled 'Advanced Global FDR Settings' is **not applicable to small**  
846 **molecule workflows**. It is currently for global False Discovery Rates across submitted  
847 files in proteomics datasets.

848

849 8) In the 'Workflow submission' section, enter an email address at which you will  
850 receive notifications when workflow jobs are completed.

851

852 9) **!CRITICAL STEP** *Making your dataset public: this is not automatic and must be*  
853 *done explicitly after submitting data and generating a dataset MSV accession number.*

854

855 Once a dataset is submitted to MassIVE, it will have an MSV accession number, and will  
856 be a private dataset in the repository, accessible only to the submitter through their  
857 personal user interface or via a user approved password protected link (e.g. perhaps  
858 during a review for publications). To make a dataset public, first select the 'Jobs' tab of the  
859 user workspace portal (**Box 1**) to find the dataset. In the list of all job submissions,  
860 MassIVE dataset submissions will appear as 'MASSIVE-COMPLETE' workflows. Click on  
861 'DONE' next to the MassIVE dataset to be made public and choose 'Make Dataset  
862 Public'. On the MassIVE website, to enable immediate use of the MassIVE dataset  
863 for GNPS workflows click on the „Convert Spectra“ tab. This converts the uploaded  
864 files to .mzML in a new folder called „ccms peak“. Otherwise, the uploaded data  
865 will be queued for this conversion and will not be immediately available.

866

867 The dataset [MSV000083437](#) has been made public, as illustrated in **SI Fig. 3**; this feature  
868 enables any reader to interact with the data and follow along with this workflow.

869

870

871 **BOX 2. The Importance of making your GNPS-MassIVE data public.**

872 Many GNPS users do not realize that when they have a dataset with MSV accession  
873 number their data is not yet public and thus remains in their private space, in accordance  
874 with GNPS-MassIVE philosophy that the data depositor should define how much and when  
875 they want to share their data in the public domain. Alternatively, upon submission, users  
876 can choose to make a dataset entirely available or 'public' to the GNPS community for  
877 browsing, commenting, subscribing, and/or downloading. This not only promotes  
878 robustness and reproducibility in MS data analysis, but also provides the user with access  
879 to the knowledge of the entire community. Indeed, the utility of GNPS for all users  
880 increases as more data becomes public, and the information and knowledge gained by  
881 any one user from this free service to the community derives from contributions made by  
882 the rest of the GNPS community. Thus, if you are a GNPS user benefiting from community  
883 contributions, by making your datasets public (and contributing network annotations,  
884 section 3.5), you are giving back to the community. It is encouraged that all users make  
885 their data public as early as possible, which provides the depositor with access to  
886 advanced features that are not available for private datasets. These features include being  
887 able to subscribe to the dataset, find related datasets, share datasets with collaborators,  
888 access living data, and utilize emerging features such as Mass Spectrometry Search Tool  
889 or MASST (the equivalent of BLAST for small molecules<sup>119</sup>). It is expected that features  
890 will continue to be developed further, thereby continually increasing the value for the end  
891 user, of both their own and other public datasets.

892

893 **3.3 Molecular networking in GNPS (SI Fig. 4) - Few minutes to several hours/days**  
894 **(depending on dataset size, user expertise)**

895

896 Once MS data files are uploaded as datasets in GNPS-MassIVE, they are available to use  
897 for analysis workflows within GNPS. Here we highlight how to execute the molecular  
898 networking workflow. A dataset can be recalled from either private or public domains in



899 MassIVE for networking analysis. Once data files have been added, they will be populated  
900 in the 'Basic Options' section of the workflow selection. The user must then input a number  
901 of parameters before running the GNPS job in both the 'Basic Options' section and in a  
902 number of 'Advanced Options' sections. The advanced parameters are dependent on  
903 analysis platform, experimental setup and conditions for acquisition of mass spectra, and  
904 will require the user to understand their ionization methods, fragmentation conditions and  
905 energies, mobile and stationary phases, and the fragmentation behavior of molecules of  
906 interest. Suggested settings for a variety of platforms are provided in the experimental  
907 section (Equipment Setup, Mass Spectrometry). A GNPS job will take approximately 10  
908 min for small datasets (up to 4 LC/MS files), 1 hr for medium datasets (5 to 400 LC/MS  
909 files), and several hrs (to days) for larger datasets (400+ LC/MS files).

910  
911

## 912 **Molecular networking workflow**

913

- 914 1) Log in to GNPS (refer to section 3.2.1 for information about how to set up account).  
915 The GNPS website banner contains tabs to navigate the platform, including tabs  
916 to navigate to MassIVE datasets, help Documentation and Forum, along with  
917 Contact information (**SI Fig. 1, Box 1**).
- 918 2) Upload desired dataset(s) to MassIVE (section 3.2.2). This step can be skipped if  
919 importing existing data files from MassIVE. Readers following the tutorial can omit  
920 this step because the GNPS-MassIVE dataset [MSV000083437](#) already exists.
- 921 3) From the GNPS splash screen (home page), start a molecular networking job by  
922 clicking the 'Create Molecular Network' button (**SI Fig. 4a**). This will bring up the  
923 main workflow input page which has a number of fillable fields to complete under  
924 each of ten sections (**SI Fig. 4b**).
- 925 4) In the 'Networking Parameter Presets' section, one of three options may be  
926 selected to set the networking parameters to approximately appropriate values  
927 depending on the size of your dataset. Clicking on one of these three options will  
928 open a workflow input form in a new tab. The default workflow settings are for  
929 'medium data'. 'Small data' refers to a dataset of up to 4 LC-MS files, 'medium data'  
930 corresponds to datasets of 5 to 400 LC-MS files, and 'large data' is applicable to  
931 datasets of more than 400 LC-MS files (e.g. [MSV000083437](#) is a medium dataset  
932 with 113 files in total). Since readers following the tutorial on the dataset  
933 [MSV000083437](#) are guided through selection of parameters, no Parameter Preset  
934 should be chosen for this example.
- 935 5) In the 'Workflow Selection' section, enter a descriptive name for the job into the  
936 'Title' field to facilitate retrieval of the workflow upon its completion. Readers  
937 following the tutorial can type 'GF/SPF Mouse Duodenum Example' in the 'Title'  
938 field (**SI Fig. 4c**).
- 939 6) Under 'Basic Options', the user will input the LC-MS files for the molecular  
940 networking workflow by choosing the 'Select Input Files' tab next to the 'Spectrum  
941 Files (Required)' field. A pop-up window with three tabs will appear: 'Select Input  
942 Files', 'Upload Files', 'Share Files' (**SI Fig. 4d**). If you are interested in analyzing  
943 multiple datasets together, you will have to repeat the above procedure with the  
944 other MSV numbers to import them into your user space.

945

946 For readers following the dataset [MSV000083437](#) tutorial, files can be imported by  
947 selecting the 'Share Files' tab. In the 'Share Files' window enter the MassIVE  
948 accession number for the dataset ([MSV000083437](#)) in the 'Import Data Share' box  
949 (**SI Fig. 4e**). After clicking 'Import', the dataset will appear in your GNPS user  
950 workspace and files can be selected for the GNPS networking workflow under the  
951 'Select Input Files' tab as described below.

- 952 7) For inputting mass spectrometry files already in your user workspace choose the  
953 'Select Input Files' tab (**SI Fig. 4f**). From the list of datasets towards the lower left  
954 of the window, select all of the files you want to analyze by clicking on individual  
955 files or an entire folder. For readers following the tutorial, GF1, GF2, GF3, GF4,  
956 SPF1, SPF2, SPF3, and SPF4 should be selected from the folder labeled 'peak'.  
957 8) Next click on the 'Spectrum Files G1' button (top of left-hand column list, with green  
958 arrow) to mark this folder / files for analysis. Your selection(s) should appear in the  
959 'Selected Spectrum Files G1' folder in the right-hand column of the window. For  
960 readers following the tutorial, folders containing data for GF1, GF2, GF3, GF4,  
961 SPF1, SPF2, SPF3, and SPF4 should now be under 'Selected Spectrum Files G1'  
962 (**SI Fig. 4g**).
- 963 9) Load the associated metadata file (see **Box 3** for format) separately into the  
964 'Selected Metadata File' folder. To do this, select the file from your workspace list  
965 (often within a MassIVE dataset in the folder labeled as 'other'), click on the  
966 'Metadata File' tab with the green arrow, and check that the file appears in the right-  
967 hand 'Selected Metadata File' folder. For readers following the tutorial,  
968 '3DMouse\_duodenum\_metadata.txt' can be selected from the folder labeled  
969 'other' (**SI Fig. 4h**).
- 970 10) Once files have been selected, the popup window can be closed by clicking on  
971 'Finish Selection'. Datasets from both your private workspace and the public  
972 domain can be recalled using either strategy. For readers following the tutorial, the  
973 final data input is shown in **SI Fig. 4i**.
- 974 11) In the 'Basic Options' section, fill in the 'Precursor Ion Mass Tolerance' (PIMT) and  
975 'Fragment Ion Mass Tolerance' (FIMT) fields taking into consideration the  
976 instrument resolution and calibration, as well as the acquisition parameters and the  
977 targeted/anticipated molecular masses (see definitions and **Table 2** below). The  
978 default is  $\pm 2.0$  Da for PIMT and  $\pm 0.5$  Da for FIMT because the reference libraries  
979 also contain spectra from low resolution instruments (e.g. ion traps of QqQ). These  
980 can be adjusted to any appropriate value. For high resolution instruments the  
981 values commonly used are  $\pm 0.01$  Da (Orbitrap) and  $\pm 0.02$  Da (qTOF) for both  
982 PIMT and FIMT.

983

984 For readers following the tutorial example, data were acquired on Bruker MaXis  
985 qTOF instrument using  $\pm 0.02$  Da. The 0.02 Da value translates into a maximum  
986 error of 40 ppm at  $m/z$  500, 20 ppm at  $m/z$  1000 for the precursor ion, and 13 ppm  
987 at  $m/z$  1500, which is consistent with the typical  $m/z$  range for small molecules.  
988 (Note that peptidic small molecules may be 2000 Da or more, although multiply  
989 charged, and thus PIMT and FIMT values of 0.03 Da should be used.) Therefore,  
990 readers should use  $\pm 0.02$  Da for both PIMT and FIMT for the example dataset (**SI**  
991 **Fig. 3c**).

992

993 **CRITICAL NOTE:** The default parameters recommended above for high resolution  
 994 mass spectrometers will not result in comprehensive searches of the spectral  
 995 libraries generated on low resolution mass spectrometers, such as ReSpect<sup>78</sup>,  
 996 large portions of MassBanks<sup>74</sup>, GNPS community contributed; a significant portion  
 997 of spectra that were annotated by matching to the NIST Mass Spectral Library with  
 998 Search Program Data Version: NIST v17 ([https://www.nist.gov/srd/nist-standard-  
 999 reference-database-1a-v17](https://www.nist.gov/srd/nist-standard-reference-database-1a-v17)) are also low resolution. In addition the natural  
 1000 products community contributes annotated spectra that may be high or low  
 1001 resolution, from a range of different spectrometers.  
 1002 **!CAUTION!** Though using low resolution parameters may increase the number of  
 1003 annotations, it will also increase the number of false positive annotations.  
 1004

1005 PIMT: This parameter is used for MS-Cluster<sup>10,12</sup> and spectral library searching, and the  
 1006 value influences the clustering of nearly identical MS<sup>2</sup> spectra via MS-Cluster.

1007 FIMT: For every group of MS<sup>2</sup> spectra being considered for clustering (consensus  
 1008 spectrum creation), this value specifies how much fragment ions can be shifted from their  
 1009 expected *m/z* values.

1010  
 1011  
 1012

**Table 3.** Absolute mass differences (Da) and associated mass error (parts-per-million, ppm) for illustrative *m/z* values

	<b>2.0 Da</b>	<b>0.5 Da</b>	<b>0.1 Da</b>	<b>0.05 Da</b>	<b>0.03 Da</b>	<b>0.025 Da</b>	<b>0.02 Da</b>	<b>0.0175 Da</b>	<b>0.015 Da</b>	<b>0.01 Da</b>	<b>0.0075 Da</b>
<b><i>m/z</i> 200</b>	10000 ppm	2500 ppm	500 ppm	250 ppm	150 ppm	250 ppm	<b>100 ppm</b>	87.5 ppm	<b>75 ppm</b>	<b>50 ppm</b>	<b>37.5 ppm</b>
<b><i>m/z</i> 500</b>	4000 ppm	1000 ppm	200 ppm	100 ppm	60 ppm	49 ppm	<b>40 ppm</b>	35 ppm	<b>29 ppm</b>	<b>20 ppm</b>	<b>15 ppm</b>
<b><i>m/z</i> 1000</b>	2000 ppm	500 ppm	100 ppm	50 ppm	30 ppm	25 ppm	<b>20 ppm</b>	17.5 ppm	<b>15 ppm</b>	<b>10 ppm</b>	7.5 ppm
<b><i>m/z</i> 1500</b>	1333 ppm	333 ppm	66 ppm	33 ppm	20 ppm	16 ppm	<b>13 ppm</b>	11.6 ppm	<b>10 ppm</b>	6.6 ppm	5.0 ppm
<b><i>m/z</i> 2000</b>	1000 pm	250 pm	50 ppm	25 ppm	15 ppm	12.5 ppm	<b>10 ppm</b>	8.75 ppm	7.4 ppm	5.0 ppm	3.75 ppm

1013  
 1014  
 1015  
 1016  
 1017  
 1018

**For advanced users:**

12) The user should complete the remaining fillable fields in 'Advanced Network Options', 'Advanced Library Search Options', and 'Advanced Filtering Options' according to their experimental design. Recommendations and values used for the example dataset are provided in Table 4 below and **SI Fig. 4j**.

- 1019 13) Use the default parameters for 'Advanced GNPS Repository Search Options',  
 1020 'Advanced Annotation Options', and 'Advanced Output Options'. The option  
 1021 'Create Cluster Buckets and BioM/PCoA Plots Output' must be enabled in the  
 1022 'Advanced Output Option' to generate bucket tables and PCoA plots from the  
 1023 'Export' and 'Advanced Views' options on the job status page (**SI Fig. 4j**).
- 1024 14) Finally, under 'Workflow Submission', the user should enter an email address to  
 1025 receive notifications when workflow jobs are completed. Readers following the  
 1026 tutorial should do this to receive notification when the example job is completed.
- 1027 15) Click 'Submit' to begin the job. The molecular networking job for the example  
 1028 dataset ([MSV000083437](#)) should take about 20 minutes.

1029  
 1030 **Table 4.** Parameters for Molecular Networking in GNPS

<b>Advanced Network Options</b>		
<b>Fillable Field</b>	<b>Definition</b>	<b>Recommended User Input</b>
Min Pairs Cos	Minimum cosine score required for an edge to be formed between nodes	Most commonly set to 0.7 when a minimum of 6 ions are matched. When fewer ions are used, it is better to be more stringent and increase this value (e.g. 0.8) but when more ions are required, one can relax this value (e.g. 0.6) <sup>120</sup> (Use 0.7 for example MSV000083437)
Minimum Matched Fragment Ions	Minimum number of common fragments that must be matched by two nodes for an edge to be formed	<i>Highly</i> dependent on the experiment – While 6 is listed as default, a lower value could be used if the user wants to be less restrictive or if the sample largely contains molecules with a small number of fragment ions. The maximum number of significant annotations are found when this value is set to 4 or 5 <sup>120</sup> . (Use 4 for example MSV000083437)
Network TopK	Maximum number of neighbor nodes for one single node. The edges between two nodes are kept only if both nodes are within each other's TopK most similar nodes. If this value is set at 10, a single node may be connected to up to 10 other nodes.	Default is set to 10. Adjusting this value enables the network to be more or less stringent. Keeping this value low makes very large networks (many nodes) much easier to visualize.  (Use 10 for example MSV000083437)
Minimum Cluster Size	Minimum number of identical MS <sup>2</sup> spectra that are merged by MS-Cluster	This is a very important parameter as it is a very good quality of spectra filter. If this is set to 1 then each MS <sup>2</sup> spectrum is compared to

	for the consensus spectrum to be represented as a node	all other MS <sup>2</sup> spectra, including MS <sup>2</sup> spectra of noise thus increasing the computational time and exploding the final molecular network. By requiring more identical spectra to be merged (clustered) before considering the MS <sup>2</sup> spectral alignments it will ensure that only reproducible and higher quality data is used in the final molecular network. The default is set to two but if it is a very large dataset (hundreds to thousands of files) one may use 5 or more while for smaller datasets (e.g. 1 or 2 files) it may be set to 1 or 2. (Use 4 for example MSV000083437)
Run MSCluster	Clusters MS <sup>2</sup> spectra and creates consensus MS <sup>2</sup> spectra using the specified mass tolerance settings	Set to 'yes' for classical molecular networking (Set to 'yes' for example MSV000083437)
Maximum Connected Component Size (Beta)	Maximum number of nodes that can be connected in a single component (molecular family) of a molecular network. This process iteratively breaks up large 'hairball' networks (of false positives) by removing the lowest scoring alignments (by cosine score) first until the resulting pieces fall below the maximum size.	Default setting is 100 – this value can be set to 0 to allow for an unlimited number of nodes or a higher setting can be used for larger datasets or for datasets containing many structurally-related molecules. (Use 100 for example MSV000083437)
Metadata File (= sample information file)	File added to the analysis that describes the experimental setup and details to allow for better downstream data visualization, analysis and interpretation	Add as a .txt file that follows the template and instructions available in the supporting information (Metadata file uploaded is described in step 9, section 3.3. Example metadata can be found in SI Tables 10 and 11, and a description of how to create a metadata file can be found in <b>Box 3</b> .)
Group Mapping and Attribute Mapping	Legacy version of metadata file	It is encouraged to use the metadata table instead
<b><i>Advanced Library Search Options</i></b>		
Library Search	Minimum number of shared	The default value is 6. Dependent on the aim

Min Matched Peaks	fragment ions to make a library match.	of the experiment: a lower value may yield more tenuous matches to library spectra, suitable for exploratory structure searching; a higher value, selecting for closer matches, facilitates dereplication of putative known compounds. The impact of this parameter is discussed in Scheubert et al. <sup>120</sup> (Use 4 for example MSV000083437)
Score Threshold	Minimum cosine similarity score to make a library match.	The default setting is 0.7. Dependent on the aim of the experiment: a lower value may yield more tenuous matches to library spectra, suitable for exploratory structure searching; a higher value, selecting for closer matches, facilitates dereplication of putative known compounds. (Use 0.7 for example MSV000083437)
Search Analogs	Matches query spectra against library spectra with a modification tolerant search within a specified range for mass differences. Precursor ion $m/z$ are allowed to deviate up to a user-defined maximum. Fragment ions that differ by the mass difference of the two parent ions are also considered.	Dependent on the user's preferences, selecting 'Do Search' requires more computing time but also the results are more exploratory. It allows for dereplication not only of identical molecules but also related molecules.
Maximum Analog Search Mass Difference	Maximum mass shift allowed between the query spectra and library spectra $m/z$ values to make a library match.	Use default parameter of 100 Da: may increase or decrease the value depending on properties such as anticipated molecular mass shift of related molecules in the samples. (e.g. 162 Da is a common mass shift for oligosaccharides). The larger this value the more likely spurious matches will be found.
<b>Advanced Filtering Options</b>		
Filter Below Std Dev	Applied before MS-Cluster. For each MS <sup>2</sup> spectrum, the 25% least intense fragment ions are collected and the std-dev is calculated, as well as the mean. A minimum peak intensity is calculated as mean + k * std-dev where k is user-selectable. All peaks below this threshold are deleted. By default, this filter is	<i>Using this filter is not recommended.</i> A default value of 0 should be used so that no filter is applied.

	inactive (value is set to 0).	
Minimum Peak Intensity	All fragment ions in the MS <sup>2</sup> spectrum below this raw intensity will be deleted.	This filter is infrequently used. Use a default value of 0 so that no filter is applied, especially if the raw intensities of your data are very low.
Filter Precursor Ion Window	All peaks in a +/- 17 Da around precursor ion mass are deleted. This removes the residual precursor ion, which is frequently observed in MS <sup>2</sup> spectra in the comparison of all spectra for molecular networking.	Apply filter, which is the default option.
Filter Library	Applies the above precursor ion window filter to the library as well.	Apply filter, which is the default option
Filter Peaks in 50 Da Window	Removes peaks that are not one of the top 6 most intense within a +/- 50 Da window.	This is commonly turned on. Dependent on the dataset: if samples contain a large number of low mass molecules or are complex mixtures containing compounds of low titer, this filtering should be turned off, as it may filter out relevant peaks that could be signals.

1031

1032

1033 **BOX 3: Sample information (metadata) collation and input** - Timing typically 1-2 hours  
 1034 for a small dataset; up to a few days for large complex metadata entries of large  
 1035 datasets<sup>121</sup>.

1036 The inclusion of a metadata (sample information) table is extremely valuable for  
 1037 interpreting the molecular network that is generated using the data. Although a time  
 1038 consuming step, it is also one of the most valuable steps for interpreting the final molecular  
 1039 network. The more time spent on curating sample information (metadata), the more useful  
 1040 the resulting molecular network will be. The metadata table links the MS files uploaded  
 1041 and selected for molecular networking analysis in GNPS with various attributes of the  
 1042 collated data based on the filename (such as "Filename.mzXML"). For instance, the  
 1043 metadata table provides the necessary information to visualize the "origin" of the detected  
 1044 metabolites when "origin" is one of the attributes used in the metadata table (e.g. column  
 1045 heading: ATTRIBUTE\_Origin). A metadata file can be created as follows:

1046 1) The metadata table must be provided as a text file (tab separated) and can be  
 1047 prepared in a text editor of choice (e.g. Microsoft Excel, Notepad++ for Windows,  
 1048 gedit for Linux, and TextEdit or TextWrangler for Mac OS) .

1049 a) When uploading metadata associated with a GNPS job, specifically  
 1050 formatted column headers are required. The first column header must be  
 1051 "filename" (no capitals, case-sensitive and no unusual characters such as  
 1052 @, #, !). **Important:** The filenames must be the filenames of the data (to  
 1053 be) uploaded to GNPS-MassIVE otherwise the metadata cannot be linked

1054 to the data. We recommend not to use any special characters such as @,  
1055 #, ! or spaces in any of the metadata fields.

1056 b) Each other column must begin with the phrase "ATTRIBUTE\_" before any  
1057 header description (e.g. ATTRIBUTE\_Origin)

1058 2) In order for sample information (metadata) to be incorporated into global  
1059 metaanalyses, the template provided in SI Table 10 should be utilized and labeled  
1060 "gnps\_metadata.tsv".

1061 There are a number of advantages to uploading a metadata table associated with a GNPS  
1062 job. When the network generated after data processing is subsequently opened in  
1063 Cytoscape, the nodes of sub-networks can be visualized based on their associated  
1064 metadata. This can be represented as a pie chart contained within each node. Additionally,  
1065 metadata can be used to color-code categories of samples when visualizing the MS<sup>2</sup>-  
1066 based statistics, such as principal coordinates analysis (PCoA) in browser using the  
1067 EMPEROR package<sup>122</sup> available in Qiime2<sup>123</sup>. This allows the user to quickly attribute the  
1068 molecular differences of the samples to certain characteristics found in the metadata. For  
1069 example, if two distinct groups appeared in the PCoA plot, it would then be possible to  
1070 color all samples of type one blue and all samples of type two red in order to determine if  
1071 this attribute could be responsible for the separation. However, it is important to note that  
1072 PCoA is only visual and doesn't give any statistical support; a PERMANOVA analysis  
1073 would have to be performed in order to actually test whether an attribute is responsible for  
1074 separation. Finally, data sharing is a vital part of modern science because it gives  
1075 opportunities for collaboration, wider scope analyses, and transparency promotes  
1076 reproducibility and thus scientific rigor. Without metadata attached, public data has less  
1077 value, will not be discovered as easily by others, and will not provide meaningful results  
1078 with MASST<sup>119</sup>. A metadata text-based search is being engineered in GNPS so that all  
1079 public data files with specific metadata entries may be re-analyzed together. When no  
1080 metadata is available, these public data will not be included in such searches. In short, the  
1081 visibility and value of data goes up by improving the amount of metadata that is uploaded.  
1082 Therefore, uploading metadata associated with the MS data to GNPS promotes a more  
1083 universal approach to science.

1084 3) In cases where you want to add a new/external metadata file (tab delimited text  
1085 format) into your workspace, under the 'Upload Files' tab: select the destination folder for  
1086 the upload on the left and drag the file for upload to the 'File Drag and Drop' box on the  
1087 right before following the same actions listed in this step. The online tutorial on metadata  
1088 formatting, including a template file, can be accessed at: [https://ccms-  
1089 ucsd.github.io/GNPSDocumentation/networking/#metadata](https://ccms-ucsd.github.io/GNPSDocumentation/networking/#metadata).

1090

1091 *Metadata format for 'ili'*<sup>111</sup>

1092 *For 2D or 3D molecular cartography using 'ili', metadata must contain the following*  
1093 *additional information. The spatial coordinates that dictate the spatial distribution of a*  
1094 *detected metabolite in a 2D (.PNG format) or 3D image (.STL format) must be included.*  
1095 *In addition to the column "filename", extra columns containing the following information:*  
1096 *"COORDINATE\_x", "COORDINATE\_y", "COORDINATE\_z", "COORDINATE\_radius"*  
1097 *have to be added. The x, y and z correspond to the 3D coordinates and the radius*  
1098 *corresponds to the approximate values of radii of the sampling points. An image viewer*  
1099 *can be used to estimate this value; for example, half of the difference between boundaries*  
1100 *of a sampling point in a horizontal or vertical dimension can be estimated. Additional*



1101 information related to 'ili can be obtained through  
1102 <https://github.com/MolecularCartography/ili>.

1103

1104

### 1105 **3.4 Visualization of the molecular network**

1106 To visualize molecular networks generated, the user can either (1) directly visualize their  
1107 network in the GNPS web browser for exploratory purposes, or (2) import data tables  
1108 generated for viewing in third party software, such as Cytoscape<sup>37</sup>, which is a free software  
1109 tool that enables visualization of the entire molecular network. These methods are  
1110 complementary to one another and the user should choose the preferred visualization  
1111 strategy based on their data analysis needs. The GNPS in-browser visualization tool is a  
1112 quick, simple way to begin analyzing data, particularly if the user wants to view and  
1113 compare MS<sup>2</sup> spectra within the network. However, in-browser visualization only allows  
1114 the user to view one molecular family (sub-network) at a time. For more advanced data  
1115 analysis and formatting options, the user can visualize their network offline in Cytoscape,  
1116 a program originally introduced by the systems biology community to allow visualization of  
1117 the complex relationships in biological sequence data. With Cytoscape, one can visualize  
1118 the chemical space that was detected in the mass spectrometry experiment as a molecular  
1119 network and provides a way to encode any property of the network (i.e. node label, shape,  
1120 color or size as well as edge label, thickness, etc.) with a metadata category (i.e. cohort,  
1121 cosine score, compound source). An online tutorial can be accessed at: [https://ccms-  
1122 ucسد.github.io/GNPSDocumentation/networking/#online-exploration-of-molecular-  
1123 networks](https://ccms-ucsd.github.io/GNPSDocumentation/networking/#online-exploration-of-molecular-networks).

1124

1125

#### 1126 **3.4.1 Molecular network visualization in browser**

1127 After completing the above molecular networking workflow, data analysis can be  
1128 performed directly in the GNPS web interface. The user can access the in-browser data  
1129 analysis options from the job status page (**Fig. 4**), several of which are described in **Table**  
1130 **5**.

1131

Job Status

**Workflow** METABOLOMICS-SNETS-V2  
 DONE [Clone] [Restart][Delete]

**Status**

Default Molecular Networking Results Views  
 [ View All Library Hits | View Unique Library Compounds | View All Clusters With IDs ]

Network Visualizations  
 [ View Spectral Families (In Browser Network Visualizer) | Network Summarizing Graphs ]

Methods and Citation for Manuscripts  
 [ Networking Parameters and Written Network Description ]

Export/Download Network Files  
 [ Download Clustered Spectra as MGF | Download GraphML for Cytoscape | Download Bucket Table | Download BioM For Qiime/Qiita | Download Metadata For Qiime | Download ili Data ]

Advanced Views - Global Public Dataset Matches  
 [ View Matches to All Public Datasets ]

Advanced Views - Third Party Visualization  
 [ View Emporer PCoA Plot in GNPS | View ili in GNPS ]

Advanced Views - Networking Graphs/Histograms  
 [ Nodes, MZ Histogram | Edges, MZ Delta Histogram | Edges, Score vs MZ Delta Plot | Library Search, PPM Error Histogram ]

Advanced Views - Misc Views  
 [ View Network, Node Centric | View Network Pairs | Networking Statistics | View Raw/Unclustered Spectra | View Compounds and File Occurrence ]

Advanced Views - Make Dataset Public Documentation  
 [ Make Public Dataset ]

Advanced Views - Experimental Views  
 [ Direct Cytoscape Preview/Download ]

**User** emgentry (emgentry.nc@gmail.com), UC San Diego

**Title** GF/SPF Mouse Duodenum Example

1132  
 1133  
 1134  
 1135  
 1136

**Figure 4.** GNPS Job Status Page.

**Table 5.** Data analysis options

<b>Data Analysis Option</b>	<b>Description</b>
View all library hits ( <b>SI Fig. 5a</b> )	View all spectra with reference database matches and assess the quality of the MS <sup>2</sup> match using the 'View Mirror Match' option. Readers following the tutorial example can view the mirror plot for cholic acid ( <b>SI Fig. 5a</b> ) in order to compare experimental spectra with library annotation. Readers can investigate mirror plots for other bile acids, as bile acid discovery is the focus of this example.
View unique library compounds ( <b>SI Fig. 5b</b> )	View all <i>unique</i> spectral matches to the reference database and perform side-by-side comparison between the query spectrum and reference spectrum. Readers following the tutorial can view query and reference spectra for cholic acid ( <b>SI Fig. 5b</b> ).
View all clusters with IDs ( <b>SI Fig. 5c</b> )	View all consensus MS <sup>2</sup> spectra that make up a node.
View spectral families ( <b>SI Fig. 5d</b> )	List of all spectral families (nodes that are connected to one another) and view individual sub-networks using in browser visualization
View EMPeror PCoA plot	Measures the binary Jaccard distance between samples based on presence/absence of molecular

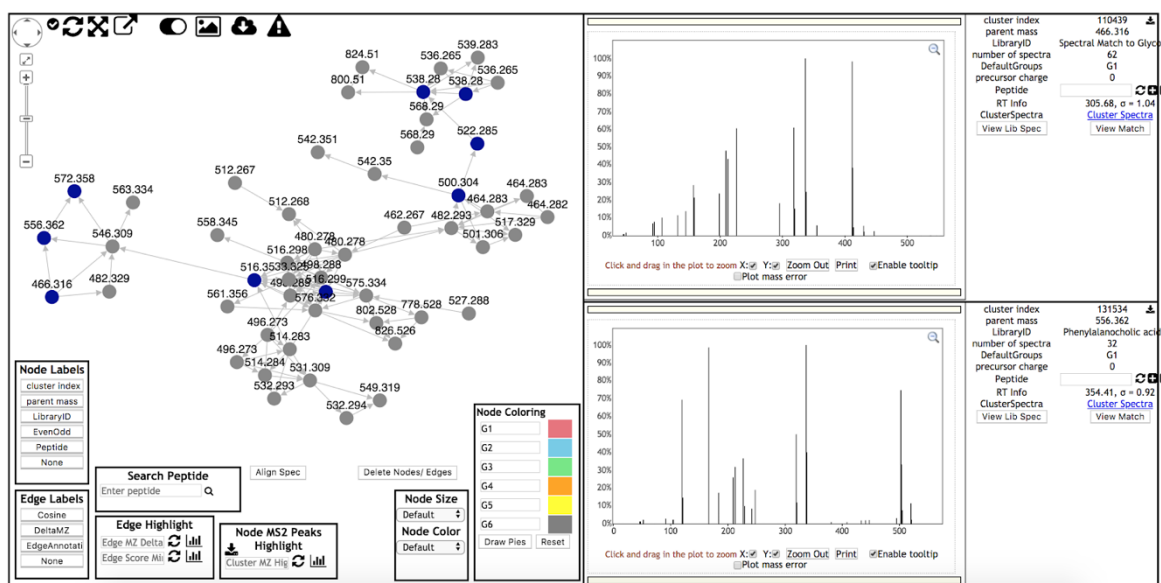
	features with associated MS <sup>2</sup> spectra as defined by the mass spectral molecular network. Interactive Principal Coordinates Analysis (PCoA) visualization is enabled through EMPeror <sup>122</sup> .
--	---

1137

1138 The “View spectral families” option lists each individual molecular family that contributes  
1139 to the entire molecular network and displays the number of MS<sup>2</sup> spectra and spectral  
1140 matches to the reference library that contribute to a given sub-network. This function also  
1141 allows users to visualize each sub-network individually in the web browser by selecting the  
1142 “Visualize network” link. Once the in browser network is displayed, the user can  
1143 immediately distinguish between nodes with library matches (blue circles) and  
1144 unannotated nodes (gray circles). Edges are represented by gray arrows that point from  
1145 the low mass spectra to the high mass spectra. Further data analysis can be performed in  
1146 this online interface as described below:

- 1147 • *Node Labels* - Nodes can be labeled by their index number given by MS-cluster,  
1148 parent mass, or library annotation name. Additionally, the node can be labeled by  
1149 a binary system to denote if the parent mass is even (1) or odd (0) to assist in  
1150 visualizing the nitrogen-rule <sup>124</sup>, or with a peptide annotation label (see Search  
1151 Peptide below). If no node label is desired, select ‘None’.
- 1152 • *Node Coloring* - This legacy feature creates pie charts to visualize mapping of  
1153 metabolites into different groups. However, this option does not use the sample  
1154 information (metadata) table and will work only if files were inputted into different  
1155 groups by the user.  
1156 !CAUTION! Note also that this is not a quantitative representation of the data  
1157 because it relies only on MS<sup>2</sup> spectral counts. Rather this feature can be used to  
1158 understand presence versus absence of compounds in specific groups.
- 1159 • *Edge Labels* - Edges connecting two nodes can be labeled with either the cosine  
1160 score or the mass difference between the parent *m/z* values (‘DeltaMZ’). If no edge  
1161 label is desired, select ‘None’.
- 1162 • *Edge Highlights* - Edges by default are represented as arrows pointing from low  
1163 mass spectra to high mass spectra, and can be colored. Users are able to enter a  
1164 mass difference (*m/z* delta) of their choice in the ‘Edge MZ Delta’ field, causing  
1165 those edges to be highlighted in red. Clicking on the graph icon next to ‘Edge MZ  
1166 Delta’ opens a new windows containing a graph that shows the distribution of all  
1167 edge *m/z* delta values in the sub-network. Selecting a peak in this ‘Network MZ  
1168 Delta Histogram’ highlights the corresponding edges in red. The same function can  
1169 be performed for ‘Edge Score Minimum’ to highlight edges that have a cosine score  
1170 greater than what the user enters.
- 1171 • *Node size/color* - The size and color of nodes can be adjusted based on spectral  
1172 counts, precursor intensity, number of files, parent mass, even/odd mass, or  
1173 precursor charge.
- 1174 • *Node MS<sup>2</sup> Peaks Highlight* - This option allows users to search the sub-network for  
1175 molecules that contain an MS<sup>2</sup> fragment of interest. To perform this query, first click  
1176 the download button within this box to pull all of the MS<sup>2</sup> spectra into the browser.  
1177 The desired *m/z* value can then be entered into the field to highlight the nodes  
1178 comprising spectra which contain the desired product ion. Alternatively, the

- 1179 histogram icon can be selected to visualize all product ions from the MS<sup>2</sup> spectra  
 1180 in the sub-network.  
 1181 • *Align Spectra* - This function enables direct comparison between the spectra of two  
 1182 connected nodes at the peak level. To perform this analysis, the user should first  
 1183 select an edge connecting two nodes, which pulls up the spectra for each node in  
 1184 the right display window. Clicking the “align spec” button overlays the spectra,  
 1185 where red peaks represent peaks of the exact same masses shared between the  
 1186 top and bottom spectra and blue peaks denote peaks matching at shifted masses.  
 1187 • *Search Peptide* - This is a function added to GNPS to support proteomic and  
 1188 peptidomic dataset analysis. If a peptide sequence is found to be associated with  
 1189 the molecular family and was found through automated peptide mining in MASSIVE  
 1190 then the amino acid sequence entered here will be searched.  
 1191



1192  
 1193  
 1194 **Figure 5.** In browser visualization of the bile acid spectral family from dataset  
 1195 MSV000083437.  
 1196

1197 **3.4.2 Assessing the quality of a library hit.** All spectral matches are putative  
 1198 annotations<sup>6</sup> until experimentally validated. Spectral matches from molecular networking  
 1199 analysis are annotations at level 2 (compounds that have been putatively annotated e.g.  
 1200 no reference standards) or 3 (compounds that can be putatively assigned to a chemical  
 1201 class based on physicochemical properties and/or spectral similarity) before validation with  
 1202 chemical standards. For level 1 annotation, the molecules would have to be isolated and  
 1203 structures elucidated or confirmed with other techniques such as NMR or X-ray analysis,  
 1204 or matching MS<sup>2</sup> and retention times, together with co-analysis with pure standards, ideally  
 1205 under more than one chromatographic condition. All non-annotated molecules in a  
 1206 molecular network are level 4 unless they are part of a molecular family containing a library  
 1207 match. Levels were defined by the 2007 Metabolomics Initiative<sup>14</sup>, and subsequently  
 1208 refined by the Compound Identification work group of the Metabolomics Society at the  
 1209 2017 annual meeting of the Metabolomics Society<sup>125</sup>. In order to judge the quality of a  
 1210 match, it is important to consider the mass accuracy of the reference spectra (resolution  
 1211 and calibration of the instrument) as compared with that of the experimental spectra. The  
 1212 sample type, experimental setup, and associated sample information (metadata) should

1213 also be taken into account when judging the accuracy of the matches. Notably,  
1214 MS<sup>2</sup> spectra typically cannot differentiate regio- or stereo-isomers and additional  
1215 experiments, including comparison with standards, are required to assign the absolute  
1216 structure.

1217 To decrease the impact of this variation all spectra, when compared, are subjected  
1218 to a square root conversion. This decreases the high intensity ions and increases the low  
1219 intensity ions. Furthermore, to address variability in data quality and source of the  
1220 reference spectra, GNPS utilizes a ranking system for submitted reference spectra, to  
1221 enable filtering of the reference library either before performing molecular networking or  
1222 afterwards, which is the default approach. Similarly the instrument that the reference data  
1223 were collected on can be considered after doing the analysis in GNPS using post-  
1224 molecular networking filtering capabilities. 'Gold' reference spectra can only be submitted  
1225 by approved users and must originate from fully characterized synthetic or purified  
1226 compounds. This is the same gold standard by which other metabolomics reference  
1227 libraries such as NIST17<sup>72</sup>, METLIN<sup>73</sup> mzCloud (<https://www.mzcloud.org/>)<sup>76</sup>, WeizMass  
1228 ([https://www.weizmann.ac.il/LS\\_CoreFacilities/weizmass-spectral-library-high-  
1229 confidence-metabolite-identification](https://www.weizmann.ac.il/LS_CoreFacilities/weizmass-spectral-library-high-confidence-metabolite-identification))<sup>126</sup> libraries are curated. Gold level spectra comprise  
1230 83% of the MS<sup>2</sup> spectra provided to GNPS as libraries. A 'silver' rating signifies that the  
1231 spectrum was submitted with an associated publication. However, GNPS also curates  
1232 crowdsourced knowledge from users in the community. All remaining reference spectra  
1233 provided by the user community receive a 'bronze' rating to denote that the annotation is  
1234 contributed by users including partial or putative annotations. The annotation within GNPS  
1235 can be made directly from the data and thus relies on the expertise of the experimentalist  
1236 and purification of the molecules is not required. This gives access to a curated reference  
1237 database that is crowdsourced and does not rely on commercially available standards. For  
1238 example, most natural products from microbes, food and plants are not commercially  
1239 available and thus the crowdsourced knowledge capture provides a resource of  
1240 information that is inaccessible any other way. The only other resource that currently  
1241 accepts putative and partial annotations is MassBank EU  
1242 (<https://massbank.eu/MassBank/>). Examples of useful but partial annotations include  
1243 modifications of molecules, such as oxidation of a molecule in which the site of oxidation  
1244 is unknown<sup>127</sup> and thus a SMILES or InChI cannot be drawn but the partial annotation  
1245 provides valuable insight to the end user. Additional partial annotations would include  
1246 adduct clusters such as sodium formate clusters or polymeric substances, including  
1247 oligosaccharides, commonly detected in mass spectrometry where a structure cannot be  
1248 drawn but is useful knowledge for the community when performing an untargeted LC-  
1249 MS/MS experiment. Users can use the above information along with the corresponding  
1250 cosine score, which takes into account the number of matching fragment ions and  
1251 differences in peak intensities, and parent mass accuracy to assess the quality of  
1252 annotation. An empirical cut-off for cosine scoring of 0.7 with 6 MS<sup>2</sup> ions matching is the  
1253 default setting in GNPS. On average this gives rise to 91% accurate annotations, and ~1%  
1254 incorrect annotations, with the remainder being attributed to possible isomers (4%) or  
1255 having not enough information by the user to judge (4%)<sup>120</sup>. However, using a target decoy-  
1256 based method to estimate confidence measures of annotations and false discovery rates  
1257 (FDR) in large scale metabolomics experiments, revealed that the annotation quality is  
1258 dataset-dependent and dependent on analysis settings such as number of ions that are  
1259 required to match. The general trend was that when few MS<sup>2</sup> ions are required to match,  
1260 a much higher cosine is required and fewer matches will be obtained at the same FDR

1261 compared to when more MS<sup>2</sup> ions are required to match the reference spectra. When more  
1262 ions are matched, the cosine score can be lowered. There is an dataset-dependent  
1263 optimum for the maximum number of spectral library matches at a specific FDR that is  
1264 typically around 4 to 6 minimum matched peaks<sup>120</sup>. Although the confidence of the spectral  
1265 matches increase when more MS<sup>2</sup> fragment peaks are required, there are fewer spectra  
1266 that have a larger number of ions, resulting in a diminished number of annotations,  
1267 especially for low MW compounds.

1268

### 1269 **3.4.3 Molecular network visualization in Cytoscape** - Timing 1-4 hours

1270 In addition to in-browser visualization, networks can be visualized using third party tools.  
1271 One popular GNPS-derived molecular network visualization tool is Cytoscape<sup>37</sup>, a  
1272 convenient software tool to use for data visualization. The steps outlined below provide  
1273 the user with a working knowledge on how to configure a network in Cytoscape. Readers  
1274 following the tutorial example can not only reproduce the same properties described in the  
1275 the steps below to generate a publishable network but also use this network to specifically  
1276 focus on the cluster containing bile acids in order to discover novel compounds.

1277

1278 There are a few options for exporting molecular networks for visualization in Cytoscape.  
1279 Once molecular networks generated from GNPS are imported into Cytoscape, a number  
1280 of simple commands can be used to make the network generated more informative,  
1281 visually appealing, and accessible (**SI Fig. 6**). [Documentation](#) on how to use Cytoscape  
1282 (versions after 3.7 release) and a [Cytoscape community forum](#) are available to assist with  
1283 troubleshooting and to learn about the latest plugins (also called Cytoscape Apps):  
1284 [https://cytoscape.org/documentation\\_users.html](https://cytoscape.org/documentation_users.html), <https://cytoscape.org/community.html>.  
1285 An online version of this tutorial is accessible at: [https://ccms-](https://ccms-ucsd.github.io/GNPSDocumentation/cytoscape/)  
1286 [ucsd.github.io/GNPSDocumentation/cytoscape/](https://ccms-ucsd.github.io/GNPSDocumentation/cytoscape/).

1287

- 1288 1. To begin using Cytoscape, download the latest version of the software from:  
1289 <https://cytoscape.org/> according to their instructions (**SI Fig. 6a**).
- 1290 2. Once Cytoscape has been downloaded, molecular networks can be imported and  
1291 visualized using two different strategies. The first option (a) will show a network  
1292 with no preset layout, while the second option (b) will show a network with default  
1293 layout settings.
  - 1294 a. In order to import data for a network with no layout present (option 1), click on  
1295 “Download GraphML for Cytoscape” in the GNPS Job status window (**SI Fig.**  
1296 **6b**). This will prompt an immediate download of a compressed folder containing  
1297 the .graphML file of interest; after uncompressing this folder using a variety of  
1298 programs, Cytoscape can be opened. The import network button (three nodes  
1299 connected by edges, **SI Fig. 6c**) in Cytoscape can be selected, permitting  
1300 selection of the .graphml file to load the network of interest.
  - 1301 b. The second option for opening a network in Cytoscape is to click on “Direct  
1302 Cytoscape Preview/Download” in the GNPS Job Status window (**SI Fig. 6d**).  
1303 This will direct the user to a new window where a pre-configured version of the  
1304 molecular network will be displayed. In this window, click on “Download  
1305 Cytoscape File” to download the file as a Cytoscape session file (.cys file) with  
1306 the visualization parameters already defined. Cytoscape can then be opened  
1307 by double clicking on the downloaded .cys file and this network will come  
1308 preloaded with GNPS default layout.

1309 c. Readers following the tutorial can use either strategy to open the completed  
1310 GNPS job run on dataset MSV000083437.

1311

1312 3. Once the molecular network has been loaded into Cytoscape, it can be customized  
1313 for viewing. By altering many properties of nodes, edges, and networks such as  
1314 colors, sizes, shapes, and labels, the default network can be transformed into a  
1315 chemically informative molecular network. Readers following the tutorial example  
1316 are guided through this process in steps 3a-3j. In the Control panel window, located  
1317 on the left side of the screen, the style and select tabs offer many options.

1318

1319 To alter a node style, click on the “Style” tab at the top of the Control Panel, then  
1320 click on the “Node” tab at the bottom of this window (**SI Fig. 6e**).

1321

1322 a. The node labels can be changed in Cytoscape by selecting the dropdown arrow  
1323 next to the “Label” tab. Readers following the tutorial example can label nodes  
1324 by selecting “precursor mass” as column and “Passthrough Mapping” for  
1325 mapping type (**SI Fig. 6f**).

1326

1327 b. Node shape can also be changed. Readers following the tutorial example can  
1328 click directly on the “Shape” symbol button and select “Ellipse” shape or change  
1329 to another desired shape (**SI Fig. 6g**). If using ellipse, the shape can be  
1330 converted into a circle by checking the box labeled ‘lock node width and height’  
(**SI Fig 6h**).

1331

1332 c. To change the node color, click on “Fill Color” dropdown. Under this column,  
1333 readers following the tutorial example can select the desired value (i.e.  
1334 “ATTRIBUTE\_host\_microbiome”) and use this to discriminate groups (i.e. germ  
1335 free vs. specific pathogen free) from one another. Readers can select “Discrete  
1336 Mapping” under the “Mapping Type” column, which allows for the selection of  
a color to be associated with each group (**SI Fig. 6i**).

1337

1338 d. Alternatively at the “Fill Color” option, the “Image/Chart 1” tab can be used to  
1339 visualize the relative ion distribution from each chosen group in the nodes as a  
1340 pie chart. Readers following the tutorial can perform this type of visualization  
1341 by clicking on the “Image/Chart 1” button, selecting the “Charts” tab, and  
1342 choosing a chart type (the pie chart is chosen in this example). The spectral  
1343 count information from groups defined in the metadata file can then be selected  
1344 from the “Available columns” to the “Selected columns” (**SI Fig. 6j**) and the user  
1345 can edit the chart color scheme using the “Options” tab. In this example, “Germ  
1346 free” and “Specific Pathogen free” can be selected and colored pink and blue,  
respectively.

1347

1348 e. To visualize the variation in the occurrence of each ion across samples (e.g.  
1349 count of 1 if not zero) as a function of the node size, go to Size option, select  
1350 “number of spectra” or “sum(precursor intensity)” as *Column* and “Continuous  
1351 Mapping” as *Mapping Type*. The opened window allows to modify the node  
1352 size in function of the node metadata column chosen. Begin by setting the value  
1353 for minimum and maximum node size value with the button *Set Min and Max*,  
1354 and then *OK*. Then move the cursor at each extremities. For readers following  
the tutorial example, set to the min size at 92 and the max at 362 (**SI Fig. 6k**).

1355

1356 f. Edge style can also be altered by clicking on the “Edge” tab at the bottom of  
the Control Panel (next to the “Node” tab) (**SI Fig. 6l**). Readers following the

- 1357 tutorial example can select this tab to make alterations in edge color and width,  
1358 in addition to other settings.
- 1359 g. To change an edge label, readers following the tutorial can click on the “Label”  
1360 dropdown arrow then select desired value. For example, mass\_difference can  
1361 be selected as “Column” in the “Passthrough Mapping mode (SI Fig. 6m).
- 1362 h. Edge width can be altered by clicking on the dropdown arrow next to “Width.”  
1363 Under the “select value” tab next to the “Column” tab, the desired value used  
1364 for scaling edges (such as cosine\_score) can be selected. At this point,  
1365 “Continuous Mapping” can be selected under “Mapping Type” (SI Fig. 6n).  
1366 Cosine\_score can be selected in the column tab and “continuous mapping” can  
1367 be chosen under mapping type to easily visualize the approximate cosine score  
1368 of all edges.
- 1369 i. The ions from experimental conditions present in the blank sample can be  
1370 subtracted from the molecular networks. In the table panel, readers following  
1371 the tutorial example can go to the column GNPSGROUP:blank, select every  
1372 rows with ion occurrence (>0), then click on the right mouse button and “select  
1373 nodes from selected rows” can be choose (SI Fig. 6o). The selected nodes  
1374 were automatically highlighted in yellow in the network. Then, do a right click  
1375 to choose in the select row “hide selected nodes and edges” (SI Fig. 6p).  
1376 However, it is possible to remove the ions from experimental conditions before  
1377 generating a molecular network by data processing<sup>128</sup>.
- 1378 j. To separate one or some specific desired network(s), press “ctrl” or “command”  
1379 (windows or MacOS, respectively) at the same time selecting the network(s)  
1380 with the mouse. Then, click on the bottom as shown in SI Fig. 6q.  
1381 Automatically, the sub-network is created. For going back to the main network,  
1382 go into the Control Panel by selecting Network, then click on the main network  
1383 bottom.
- 1384 4. At this point, readers following the tutorial example have generated a publishable  
1385 network in Cytoscape from the output of molecular networking in GNPS. This  
1386 network should look like that shown in Fig. 3. Interested readers can look more  
1387 closely at the sub-network containing key bile acids in order to practice manual  
1388 propagation of annotations throughout a sub-network (Fig. 3). Style options are  
1389 described in more detail in the Cytoscape manual:  
1390 <http://manual.cytoscape.org/en/stable/Styles.html>.

### 1391 3.5 How to propagate annotations through manual interpretation of the networks

1392 A molecular network can be very useful in propagating annotations through manual  
1393 interpretation of networks in parallel with raw MS<sup>2</sup> spectra. Manual annotation can be  
1394 performed by looking at mass differences (deltas) in the molecular network and assigning  
1395 the source of these deltas, i.e. charge retention fragmentations such as retro-Diels Alder  
1396 reactions or McLafferty rearrangements and charge migration fragmentations such as  
1397 simple inductive cleavages or  $\alpha$ - or  $\beta$ -eliminations<sup>129</sup>. The novel bile acids found in the  
1398 mouse duodenum provide an example of the utility of manual interpretation of networks  
1399 (SI Fig 7b). One can use the mass deltas between unknown nodes and neighboring library  
1400 hits to determine new structures. In the above example, three unknown nodes were  
1401 determined to be novel bile acids conjugated with phenylalanine, leucine, and tyrosine  
1402 based on their mass deltas with respect to glycocholic or glycomuricholic acid. A  
1403



1404 description of how manual propagation of annotations can be performed in the context of  
1405 the example is given below:

1406

- 1407 1) The Cytoscape's toolbar can be used to search nodes or edge metadata (e.g.,  
1408 "shared name"). Readers following the tutorial example can enter "glycocholic acid"  
1409 with the quotation marks. The node of interest at  $m/z$  466.316 that matches  
1410 glycocholic acid in the GNPS library are automatically selected and highlighted in  
1411 yellow in the network (**SI Fig. 6g**).
- 1412 2) Manually propagate annotation based on mass shifts. In **SI Fig. 7a**, glycocholic  
1413 acid connects to a node with  $m/z$  556.363. Based on the mass shift of 90.047, the  
1414 unknown node can be manually annotated as glycocholic acid conjugated with  
1415 phenylalanine. Analogously, nodes with  $m/z$  572.358 and 522.379 could be  
1416 manually annotated as glycocholic acid conjugated with tyrosine and leucine  
1417 respectively, accounting for mass shifts of 106.042 and 56.063 Da.
- 1418 3) The select function is helpful to find the annotated nodes within the network with a  
1419  $m/z$  error from 0 to 10 ppm between precursor ions. This tool is available in Control  
1420 Panel at the Select tab, and can be used to create a selection of nodes and/or  
1421 edges based on their metadata and/or network topology. Readers following the  
1422 tutorial example can click on the "+" button and choose "MZErrorPPM" as column  
1423 filter and move the cursor from 0 to 10, then click on Apply (**SI Fig. 7b**). These  
1424 nodes are automatically selected and highlighted in yellow in the network.
- 1425 4) Advanced computational tools can also be used for automated annotation  
1426 propagation, such as the Network Annotation Propagation (NAP) tool<sup>85</sup>, or manual  
1427 annotation can be performed using the results of Dereplicator<sup>82, 83</sup> and  
1428 Mass2Motifs,<sup>130</sup> which can be accessed through GNPS at  
1429 <https://gnps.ucsd.edu/ProteoSAFe/static/gnps-theoretical.jsp>.

1430

### 1431 3.6 Capturing information by adding reference spectra from your data

1432 Once an MS<sup>2</sup> spectra has been fully annotated, it can be added as a reference  
1433 spectrum to GNPS. Because the GNPS library database is crowd-sourced, users are  
1434 encouraged to submit spectral annotations because knowledge they have is captured  
1435 through these annotations of reference spectra and reusable by others. This enables the  
1436 creation of reference spectra from MS<sup>2</sup> spectra in the dataset without needing to purify the  
1437 molecule. The assumption is made that the people who collected the data are experts in  
1438 their samples and thus are in the best position to curate. Additionally, if the same user or  
1439 lab then uploads another related dataset, and it contains the same molecule, it will be  
1440 automatically annotated. Users can upload a single reference spectrum by first clicking on  
1441 "View All Clusters With IDs" in the job status page, then selecting the cluster desired for  
1442 annotation from the "ClusterIdx" column. Once the cluster is selected, the  
1443 "AnnotatetoGNPS" button can be selected. This button brings up the workflow for  
1444 annotation, where input files, sample parameters, desired annotation, advanced  
1445 annotations and library selections can be added and the job can be submitted. Users can  
1446 also add a known spectrum to the library from a file uploaded to MassIVE by selecting  
1447 "Contribute" under the "Add Your Spectrum" heading on the main page, even if molecular  
1448 networking has not been performed on this file. Additionally, if the user wishes to upload  
1449 >50 reference spectra to GNPS, a separate batch upload can be performed to streamline  
1450 the process as detailed in the online help [documentation at https://ccms-](https://ccms-ucsd.github.io/GNPSDocumentation/batchupload/)  
1451 [ucsd.github.io/GNPSDocumentation/batchupload/](https://ccms-ucsd.github.io/GNPSDocumentation/batchupload/). All annotations can be refined at a later

1452 step, and the provenance of each curation is retained within the GNPS-MassIVE  
1453 environment. For example one person may annotate that they think it is a lipid, the next  
1454 person may update and specify it is a phosphatidylcholine and the next person may refine  
1455 this to be 1-oleoyl-2-palmitoyl-phosphatidylcholine and this is all logged in the CCMS  
1456 spectral library for each MS<sup>2</sup> spectrum.

1457

### 1458 **3.7 Data sharing & reproducibility of molecular networking**

1459 GNPS users are encouraged to share both the raw mass spectrometry data and  
1460 associated molecular networking jobs that contributed to peer-reviewed publications by  
1461 providing the MassIVE accession number (e.g. MSV000083437) and a hyperlink to the  
1462 GNPS job in the methods or experimental details section of the publication. Datasets  
1463 uploaded to MassIVE ideally include all raw and peak picked mass spectrometry data and  
1464 associated sample information (metadata). GNPS records all data inputs, manipulations  
1465 and analyses of the data, providing a historical record of the data and its origins. This data  
1466 provenance promotes reproducibility and ultimately quality of the data and its annotations.

1467

#### 1468 **3.7.1. Cloning a job**

1469 Once a job's URL address is shared, any GNPS user can clone the job by following  
1470 the provided link and clicking 'clone' on the job status page (**SI Fig. 8**). Cloning a job allows  
1471 users to view all parameters and files that were used to create the existing network and  
1472 easily rerun the molecular networking job with the same (or adjusted) parameters and files.  
1473 Cloning a GNPS job is an extremely useful tool that promotes reproducibility and scientific  
1474 rigor. This is a feature many users use to submit multiple molecular networking jobs with  
1475 modified parameters. Note that if data were imported from your private user workspace  
1476 and not from within MASSIVE, other users will not have access to the mass data and  
1477 consequently will not alter the analysis in GNPS. If a job has been run in the previous V1  
1478 version of GNPS (i.e. it ran using the 'METABOLOMICS-SNETS' workflow), it can be  
1479 cloned and re-run in version 2 (V2) of GNPS by simply clicking 'Clone Job to Latest  
1480 Molecular Networking V2 Workflow' on the job status page (**SI Fig 8b**).

1481

#### 1482 **3.7.2. Accessing a dataset**

1483 If a dataset is public, users are able to download all files for reanalysis, including  
1484 raw data and the sample information table (metadata). To access a MassIVE dataset of  
1485 interest, users should select 'MassIVE Datasets' in the GNPS workspace portal (**Box 1**)  
1486 and enter the MassIVE accession number or defining keywords into the search bar. The  
1487 user can then click on the MassIVE accession number highlighted in green to link to the  
1488 'MassIVE dataset information page', and select the 'FTP Download' link to download files.  
1489 Alternatively, this link can be pasted into the quick connect box of an FTP client.

1490 In contrast, private datasets can only be viewed by the user who uploaded the data  
1491 and anyone who has a link to the job status page. The user can create a password  
1492 protected link. When downloading data from a private dataset, you will be prompted to  
1493 enter a password for that MassIVE dataset ID. If using an FTP client, you will need to enter  
1494 the MassIVE ID as the username, followed by a password. If the submitter did not specify  
1495 a password, then it should be accessible using the password 'a'.

1496

#### 1497 **3.7.3. Subscribing to a dataset and living data**

1498 Public datasets remain alive long after publication: for example, they will be  
1499 searched periodically against the ever growing annotated GNPS spectral libraries,

1500 potentially yielding new putative annotations within those datasets. Beyond new  
1501 identifications within a dataset, subscribers will receive email notifications of other datasets  
1502 that exhibit chemical similarities to the subscribed dataset. This allows for users to be  
1503 connected via their research interest to similar datasets. Updates are sent out about once  
1504 a month and only when there is new information associated with the dataset. To subscribe  
1505 to a dataset, the user should navigate to the 'MassIVE dataset information' page as  
1506 described above in section 3.7.2 and click 'Subscribe'. This feature changes the way we  
1507 interact with data. Previously, data was periodically reanalyzed by the submission of new  
1508 jobs, but in GNPS, data is automatically reanalyzed and updates are sent to the  
1509 subscribers. Therefore, data may give rise to useful results a few weeks or even a few  
1510 years later after it is uploaded or it may enable the dissemination of all the knowledge of  
1511 this dataset to all lab members or collaborators.

1512

#### 1513 **BOX 4: Feature-based molecular networking (FBMN)**

1514 The above described molecular networking analysis represents the type of  
1515 molecular networking that is most widely used currently. This workflow connects clustered  
1516 MS<sup>2</sup> spectra as nodes based on spectral similarities and makes use of MS<sup>2</sup> data only, even  
1517 for quantitation. The chromatographic dimension and MS<sup>1</sup> data are not considered in  
1518 classical molecular networking. However, in MS-based metabolomics studies, statistical  
1519 analysis is done predominantly from MS<sup>1</sup>-based peak abundances from extracted ion  
1520 chromatograms (XIC). These chromatographic peaks with a specific accurate mass-to-  
1521 charge ratio are described as features. In order to bridge this gap between MS<sup>1</sup> abundance  
1522 and MS<sup>2</sup> qualitative information, there is a workflow to link MS<sup>1</sup> intensities derived from  
1523 LC-MS features with MS<sup>2</sup> information from molecular networking<sup>131, 132</sup>. This workflow is  
1524 called feature-based molecular networking ([https://ccms-  
1525 ucsd.github.io/GNPSDocumentation/featurebasedmolecularnetworking/](https://ccms-ucsd.github.io/GNPSDocumentation/featurebasedmolecularnetworking/)) and can be  
1526 performed using open access mass spectrometry processing tools such as MZmine 2<sup>113</sup>,  
1527 XCMS<sup>79</sup>, MS-DIAL<sup>133</sup>, or OpenMS<sup>112</sup>. In this workflow, feature finding is the computational  
1528 process of selecting and identifying features in the MS<sup>1</sup> across multiple samples and must  
1529 be performed prior to generating a network. These tools allow the export of a feature table  
1530 and corresponding MS<sup>2</sup> scans for each feature, which can be submitted to feature-based  
1531 molecular networking through GNPS. Furthermore, the integration in MZmine 2 allows a  
1532 direct submission to GNPS even without being a registered GNPS user. However, by  
1533 providing the username and password, the new networking job is directly created in the  
1534 specified user space.

1535

#### 1536 **4.0 Troubleshooting**

1537 Table 6 below lists some more common scenarios or questions encountered when using  
1538 GNPS. We also recommend to check the forum link from the banner in GNPS where users  
1539 can post questions to the GNPS community.

1540

1541 **Table 6.** Troubleshooting

This protocol does not address the issues that the user faces.	Check the GNPS forum and post questions.
Job fails with the	Check that your data are in a supported file format; check that the

message 'Empty MS/MS'	submitted files are centroided and have MS <sup>2</sup> data, check that filtering criteria are not too aggressive; check that raw files are not included in the file selection.
Job fails with the message 'spectral library search exceeded memory'	This means that the spectral library search step used too much memory and had to be terminated. This is likely caused by changing the set of spectral libraries used in search (such as removing the spectra filtering). This issue can potentially be resolved by increasing the maximum cluster size value to reduce the number of searched spectra. It is not recommended to change the set of libraries included unless you are an advanced user. Please remove all libraries except for the default "speclibs" and rerun.
Network is too large to view in Cytoscape	If a dataset cannot be loaded into Cytoscape, a sub-network of interest can be opened. Alternatively, larger networks can be opened on a computer with more RAM.
I do not know how to include / exclude blanks	This is most easily addressed if the blanks are included in the metadata. Then the user can opt to visualize spectra found in blanks using discrete mapping in Cytoscape or other visualization tool.
Metadata does not sync with Cytoscape	The metadata (sample information) table must be formatted correctly. In particular, check whether the first column is named 'filename', whether all file names match exactly the files uploaded to GNPS and have '.mzXML' extensions (or other compatible file format), and whether each metadata column uses the prefix "ATTRIBUTE_" and there are no trailing spaces in any of the headings.
GNPS job fails due to improper metadata format	The metadata file must be formatted as a tab-separated .txt file.
I cannot see my file(s) after drag and drop upload to GNPS workspace	Check that the targeted folder is highlighted before dragging and dropping file.
My GNPS network is much smaller (fewer nodes) than expected	Check that you selected the mzXML peaklist files from the 'ccms_peak' folder of your MassIVE dataset for the GNPS workflow, not the mzXML files generated directly from the raw data files in the 'raw' folder. The value of the minimum cluster size can be reduced. The minimum cosine score can also be decreased to increase the number of edges in the networks.
Cannot convert Waters .raw files to .mzXML / .mzML from data acquired in the MSE mode of Waters mass	Datasets acquired on a Waters mass spectrometer using the MSE mode can currently only be converted to .mzML using the vendors UNIFI platform. Alternatively, data need to be collected in DDA or MS <sup>2</sup> mode, for which data conversion to .mzXML/.mzML is enabled through ProteoWizard.

spectrometers using ProteoWizard	
Molecules that I know are structurally similar do not appear to form a cluster	Check consensus spectra for the molecules of interest. It is possible that low abundance noisy spectra are included which results in poor consensus. For some classes of compounds that do not fragment efficiently, e.g. certain lipids, the MS <sup>2</sup> spectra are not informative enough to build meaningful network.

1542

1543

## 5.0 Anticipated Results

1544

1545

1546

1547

1548

1549

1550

1551

1552

1553

1554

1555

1556

1557

1558

1559

1560

1561

1562

1563

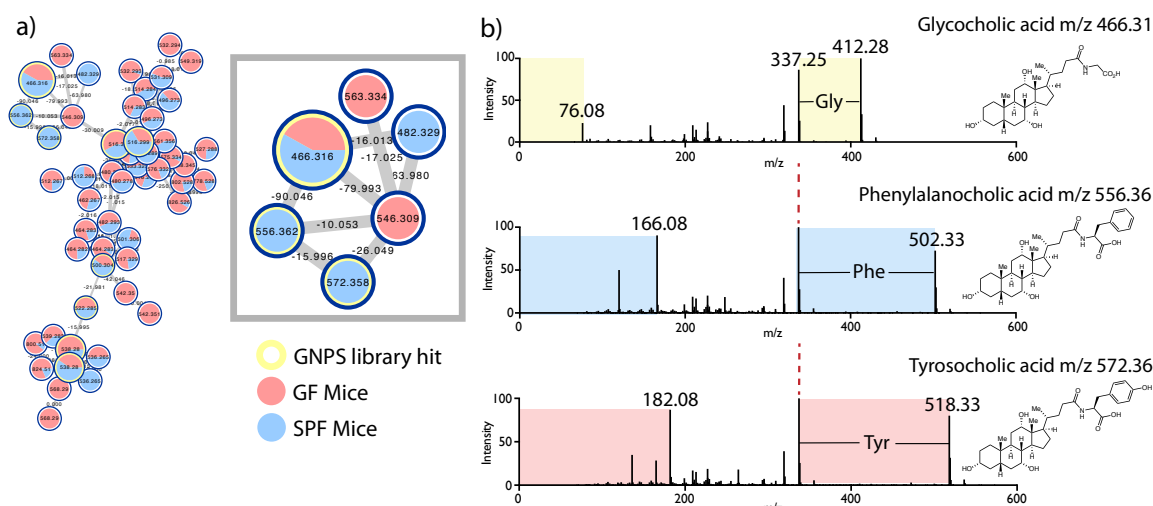
1564

1565

1566

1567

Molecular networking of LC–MS/MS data according to the protocol described herein integrates an associated sample information table (metadata file) with the latest molecular networking workflow, to yield a network (.graphml file) that may be visualized directly in GNPS or imported into Cytoscape. The tutorial example followed throughout the protocol demonstrates how contemporary GNPS molecular networking can be used to discover a new set of conjugated bile acids from the mouse gut microbiome as described in section 3.5.<sup>65</sup> The network produced from the protocol should contain a molecular family of conjugated bile acids that includes a library hit for glycocholic acid (Figure 5a). This annotation can be propagated to identify new bile acids by converting the mass differences of the edges into structural motifs. For instance, the user can identify the *m/z* 546.309 node as a sulfated cholic acid by using its mass difference of 79.993. This strategy was key in determining the structures for the new phenylalanine (*m/z* 556.362) and tyrosine (*m/z* 572.358) conjugated cholic acids. This example also showcases how manual comparison of the MS/MS spectra that make up the conjugated bile acid molecular family can also be critical for structural annotation. For example, spectra of Gly-, Phe-, and Tyr-conjugated cholic acid all contain fragment ions identical in mass to their respective amino acid conjugates (Figure 5b). Furthermore, the mass difference between the precursor ion and the common peak at *m/z* 337.25, which corresponds to amide bond cleavage, matches the exact mass of the conjugated amino acid. In addition to the conjugated bile acids, the user can also find hits for cholic acid and deoxycholic acid in the network. These compounds are present only in colonized mice, as microbes deconjugate tauro- and glyco-conjugated bile acids in the duodenum.

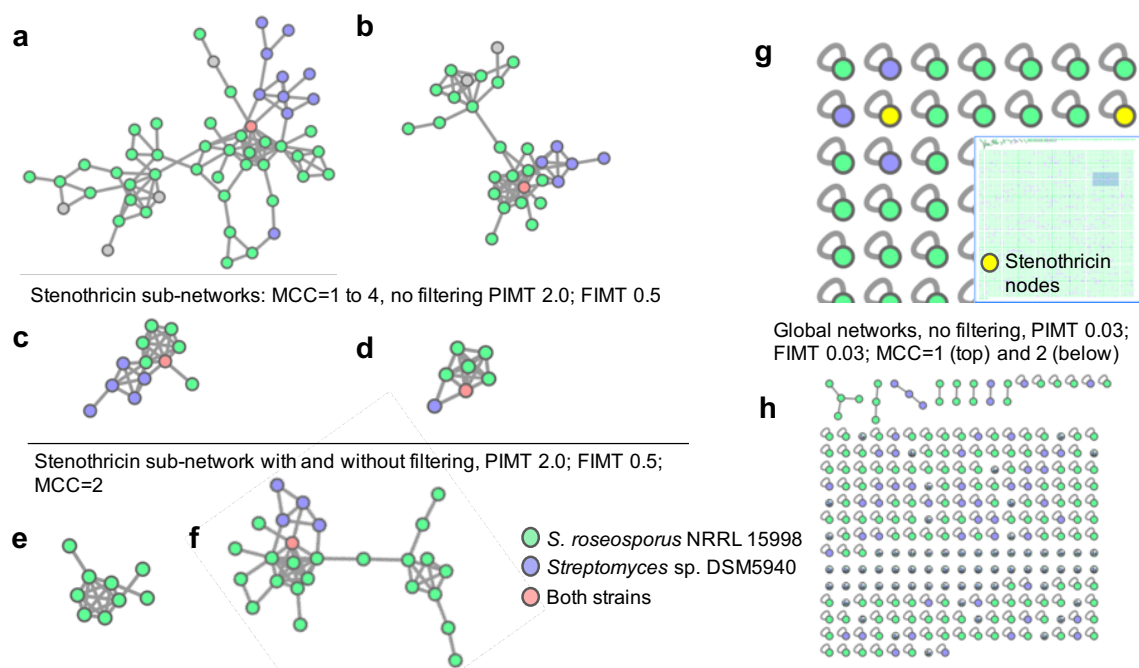


1568  
 1569 **Figure 5.** (a) The molecular family of conjugated bile acids from the duodenum of germ  
 1570 free (GF) (red) vs. specific pathogen free (SPF) (blue) mice in [MSV000083437](#)  
 1571 dataset. As shown in the inset, a library hit for glycocholic acid ( $m/z$  466.316) is present in both GF  
 1572 and SPF mice while the new phenylalanine ( $m/z$  556.362) and tyrosine ( $m/z$  572.358)  
 1573 conjugated bile acids are seen only in colonized mice. (b) Comparison of MS<sup>2</sup> spectra for  
 1574 Gly-, Phe-, and Tyr-conjugated bile acids.

1575

1576 In addition to the tutorial example, which highlights how molecular networking can be used  
 1577 for the discovery of new endogenous metabolites related to human health, two more  
 1578 examples are presented from published studies<sup>11, 50</sup>. One highlights the use of molecular  
 1579 networking in natural products discovery and the other integrates metabolomic and  
 1580 microbiome data into 3D maps. It is worth noting that the molecular networking workflow  
 1581 in GNPS continues to be updated and additional reference library entries are continually  
 1582 added by the GNPS community, which may result in some new network annotations since  
 1583 the original publication. The current reference libraries used (curated in speclibs,  
 1584 December 2018) are listed in the supporting information (SI Table 11). To illustrate the  
 1585 utility of GNPS in revealing the extent of suites of related natural products, the discovery  
 1586 of new stenothricins-GNPS 1-5 from *Streptomyces* strains reported in Wang et al.<sup>11</sup> is  
 1587 revisited here. The dataset MSV000083381 comprises MS<sup>2</sup> data for *n*-butanol and  
 1588 methanol extracts from each of *Streptomyces* sp. DSM5940 and *S. roseosporus* NRRL  
 1589 15998 cultures grown on solid agar, together with a metadata table that links each of the  
 1590 four MS<sup>2</sup> data files with the originating *Streptomyces* strain. In reproducing the observation  
 1591 of a distinct sub-network comprising the MS<sup>2</sup> data from *Streptomyces* sp. DSM5940  
 1592 connected to known *S. roseosporus* stenothricin analogs, we highlight the effect of  
 1593 minimum consensus cluster size, PIMT and FIMT settings, and advanced filtering options  
 1594 (**Fig. 6**). Importantly, the choice of low resolution settings for PIMT (2.0) and FIMT (0.5) to  
 1595 facilitate library searching enables annotation of multiple stenothricin analogues in an  
 1596 expansive sub-network, which is otherwise lost with more stringent mass tolerance  
 1597 settings of 0.03. Minimum consensus cluster size also has a pronounced effect on the  
 1598 range of stenothricin analogues detected. As is common for many natural product  
 1599 molecular families, a few major stenothricin analogues are likely accompanied by  
 1600 numerous minor stenothricins, for which the MS<sup>2</sup> spectra generated readily fall below the  
 1601 threshold for representation as a node. The distinct clustering of stenothricins from  
 1602 *Streptomyces* sp. DSM5940 in **Fig. 6a** is because the parent ion  $m/z$  values for these

1603 nodes are 41 Da less than the corresponding values for the known *S. roseosporus*  
 1604 stenothricin compounds, consistent with the substitution of serine for lysine in stenothricin-  
 1605 GNPS 1-5<sup>11</sup>.

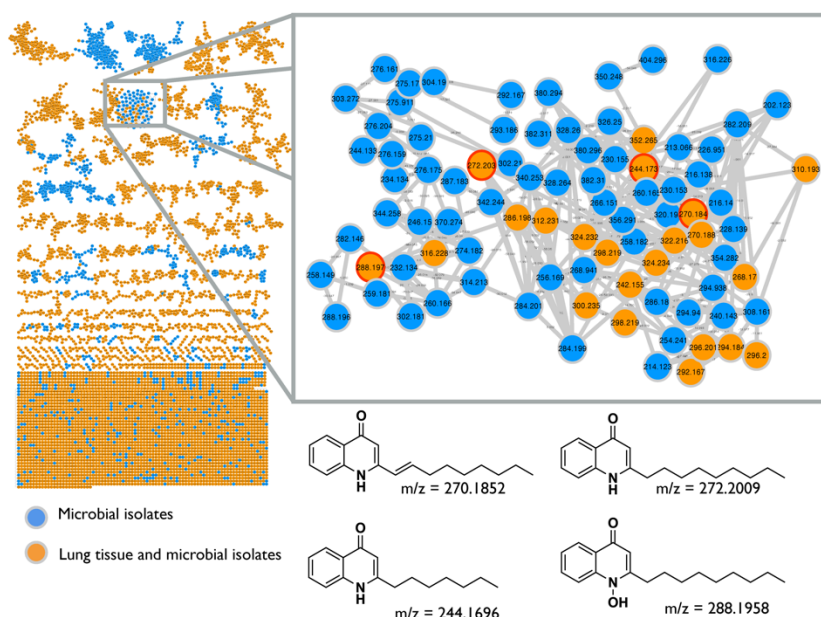


1606  
 1607

1608 **Figure 6.** Networking of the stenothricin natural product molecular family  
 1609 ([MSV000083381](#)) detected in *Streptomyces* sp. DSM5940 (purple nodes), *S. roseosporus*  
 1610 NRRL 15998 (green nodes) or both strains (yellow nodes). Variation in number of nodes  
 1611 and spectra with 'Minimum Cluster Size' (MCS) yields sub-networks (a) MCC=1, 52 nodes,  
 1612 169 spectra, (b) MCC=2, 29 nodes, 144 spectra, (c) MCC=3, 12 nodes, 89 spectra, (d)  
 1613 MCC=4, 7 nodes, 73 spectra (no filtering). Selecting advanced filtering options results in  
 1614 (e) 9 nodes, compared to (f) 26 nodes. High resolution settings for PIMT (0.03) and FIMT  
 1615 (0.03) reduce stenothricin annotations with (g) MCC = 1 providing two stenothricin nodes  
 1616 of 7642 total, and (h) MCC = 2 giving no stenothricin annotations and only 192 nodes.  
 1617 Parent ion mass tolerance = PIMT and fragment ion mass tolerance = FIMT.

1618

1619 To further illustrate that molecular networking in GNPS can be used for a diverse range of  
 1620 applications, we highlight that molecular networking can be used to visualize quinolones  
 1621 produced by *Pseudomonas* isolated from a patient lung<sup>50</sup>. **Fig. 7** reproduces the previous  
 1622 analysis ([MSV000083359](#)), where the orange nodes represent quinolones detected in both  
 1623 lung tissue extracts and cultured microbial isolates, while cyan nodes represent those only  
 1624 detected in cultured microbial isolates.



1625

1626

1627

1628 **Figure 7.** Molecular family (a sub-network) of quinolones detected in lung tissue extracts

1629 (orange nodes) and cultured *Pseudomonas* isolates (cyan nodes), created from MASSIVE

1630 dataset [MSV000083359](https://massive.ucsf.edu/MSV000083359). 2-heptyl-4-quinolone (HHQ), 2-nonyl-4-quinolone (NHQ) and its

1631 unsaturated derivative (NHQ-C9:1 db), and 2-nonyl-4-quinolone-N-oxide (NQNO) were

1632 found in lung tissue, and are highlighted by a red node border.

1633

1634 With a network in hand, there are a number of data analysis tools and experimental

1635 validation steps that may be performed. As discussed in section 3.4.2, to legitimize a library

1636 annotation beyond inspecting mirror plots, the user should verify molecular formula and

1637 identify associated adducts using MS<sup>1</sup> data. Additionally, rationalization based on

1638 biological source is recommended. Ideally, an annotation is authenticated by comparison

1639 with a known standard compound or isolation and full characterization. In the example

1640 followed throughout the protocol, the molecular structures of the new conjugated bile acids

1641 from the mouse duodenum were confirmed by comparison with synthetic standards. For

1642 more complex structures such as those in the stenothricin example<sup>11</sup> (Figure 6), the most

1643 abundant analog, stenothricin-GNPS 2, was purified for acquisition. The structure was

1644 assigned from 1D and 2D NMR data, Marfey's analysis<sup>134</sup>, and manual comparison of the

1645 MS<sup>2</sup> spectra with MS<sup>2</sup> spectra for previously reported stenothricin D. Genome mining

1646 further supported the conclusion that the -41 Da mass shift observed for stenothricin-

1647 GNPS 1-5 is due to a Lys to Ser substitution. For nodes that are not annotated, the *in*

1648 *silico* Dereplicator may predict peptidic natural products, while NAP (Network Annotation

1649 Propagation) can use annotated nodes to predict related metabolites. Molecular formulas

1650 may be generated using additional tools, one of which is SIRIUS<sup>108</sup>. This software uses

1651 MS<sup>2</sup> features to arrive at the best molecular formula for the precursor MS<sup>1</sup> ion, and works

1652 best for smaller molecules (<600 Da).

1653 In the example of the human lung colonized by *Pseudomonas* bacteria (Figure 7)<sup>50</sup>, the

1654 authors use spatial mapping to visualize annotated molecules on an exploded lung, and

1655 then correlate the distribution of molecules to microbiome maps generated from 16S rRNA

1656 gene amplicon sequencing. This study shows how molecular networking can be used to



1657 elucidate spatial variation in chemical profile and how this can be correlated with microbial  
1658 makeup using 3D maps. Statistical analyses of microbiome sequence data were  
1659 performed in QIIME2; a number of additional statistical tools as well. Ongoing  
1660 developments in GNPS include the integration of some of these statistical analysis tools  
1661 into GNPS. Ultimately, it is envisioned that streamlined integration of pre- and post-  
1662 networking tools with the GNPS platform will facilitate both creation and mining of  
1663 molecular networks.

1664

1665

#### 1666 **Acknowledgements:**

1667 National Research System (SNI) of SENACYT Panama funded CABP, CMH, JL-B, MG;  
1668 Gordon and Betty Moore Foundation (PD, NB, KLM), National Institutes of Health  
1669 (GM122016-01: KLM), National Science Foundation (DEB1354944: RMT); AKJ  
1670 recognizes the American Society for Mass Spectrometry 2018 Postdoctoral Career  
1671 Development Award. DP was supported through Deutsche Forschungsgemeinschaft  
1672 (DFG) with grant PE 2600/1. R03 CA211211 (PD) on reuse of metabolomics data and  
1673 P41 GM103484 (PD, NB) Center for Computational Mass Spectrometry as well as  
1674 Instrument support through NIH S10RR029121 (PD).

1675

1676

1677

#### 1678 **References:**

- 1679 1. Watrous, J. et al. Mass spectral molecular networking of living microbial colonies.  
1680 *Proc Natl Acad Sci U S A* **109**, E1743-1752 (2012).
- 1681 2. Traxler, M.F. & Kolter, R. A massively spectacular view of the chemical lives of  
1682 microbes. *Proc Natl Acad Sci U S A* **109**, 10128-10129 (2012).
- 1683 3. Ramos, A.E.F., Evanno, L., Poupon, E., Champy, P. & Beniddir, M.A. Natural  
1684 products targeting strategies involving molecular networking: different manners,  
1685 one goal. *Natural Product Reports Advance article* (2019).
- 1686 4. Teta, R. et al. A joint molecular networking study of a Smenospongia sponge and  
1687 a cyanobacterial bloom revealed new antiproliferative chlorinated polyketides. *Org.*  
1688 *Chem. Front.* **6**, 1762-1774 (2019).
- 1689 5. Kalinski, J.J. et al. Molecular Networking Reveals Two Distinct Chemotypes in  
1690 Pyrroloiminoquinone-Producing Tsitsikamma favus Sponges. *Mar Drugs* **17**  
1691 (2019).
- 1692 6. Raheem, D.J., Tawfike, A.F., Abdelmohsen, U.R., Edrada-Ebel, R. & Fitzsimmons-  
1693 Thoss, V. Application of metabolomics and molecular networking in investigating  
1694 the chemical profile and antitrypanosomal activity of British bluebells  
1695 (*Hyacinthoides non-scripta*). *Sci Rep* **9**, 2547 (2019).
- 1696 7. Trautman, E.P., Healy, A.R., Shine, E.E., Herzon, S.B. & Crawford, J.M. Domain-  
1697 Targeted Metabolomics Delineates the Heterocycle Assembly Steps of Colibactin  
1698 Biosynthesis. *J Am Chem Soc* **139**, 4195-4201 (2017).
- 1699 8. Vizcaino, M.I., Engel, P., Trautman, E. & Crawford, J.M. Comparative  
1700 metabolomics and structural characterizations illuminate colibactin pathway-  
1701 dependent small molecules. *J Am Chem Soc* **136**, 9244-9247 (2014).
- 1702 9. Nguyen, D.D. et al. Indexing the *Pseudomonas* specialized metabolome enabled  
1703 the discovery of poaeamide B and the bananamides. *Nature Microbiology* **2**, 16197  
1704 (2016).
- 1705 10. Frank, A.M. et al. Clustering millions of tandem mass spectra. *J Proteome Res* **7**,  
1706 113-122 (2008).

- 1707 11. Wang, M. et al. Sharing and community curation of mass spectrometry data with  
1708 Global Natural Products Social Molecular Networking. *Nat Biotechnol* **34**, 828-837  
1709 (2016).
- 1710 12. Frank, A.M. et al. Spectral archives: extending spectral libraries to analyze both  
1711 identified and unidentified spectra. *Nat Methods* **8**, 587-591 (2011).
- 1712 13. De Vijlder, T. et al. A tutorial in small molecule identification via electrospray  
1713 ionization-mass spectrometry: The practical art of structural elucidation. *Mass*  
1714 *Spectrom Rev* **37**, 607-629 (2018).
- 1715 14. Sumner, L.W. et al. Proposed minimum reporting standards for chemical analysis  
1716 Chemical Analysis Working Group (CAWG) Metabolomics Standards Initiative  
1717 (MSI). *Metabolomics* **3**, 211-221 (2007).
- 1718 15. Su, G., Morris, J.H., Demchak, B. & Bader, G.D. Biological network exploration with  
1719 Cytoscape 3. *Curr Protoc Bioinformatics* **47**, 8 13 11-24 (2014).
- 1720 16. Smoot, M.E., Ono, K., Ruscheinski, J., Wang, P.L. & Ideker, T. Cytoscape 2.8: new  
1721 features for data integration and network visualization. *Bioinformatics* **27**, 431-432  
1722 (2011).
- 1723 17. Sandhu, C. et al. Evaluation of Data-Dependent versus Targeted Shotgun  
1724 Proteomic Approaches for Monitoring Transcription Factor Expression in Breast  
1725 Cancer. *Journal of Proteome Research* **7**, 1529-1541 (2008).
- 1726 18. Hubert, J., Nuzillard, J.-M. & Renault, J.-H. Dereplication strategies in natural  
1727 product research: How many tools and methodologies behind the same concept?  
1728 **16**, 55-95 (2017).
- 1729 19. Rochat, B. Proposed Confidence Scale and ID Score in the Identification of Known-  
1730 Unknown Compounds Using High Resolution MS Data. *J Am Soc Mass Spectrom*  
1731 **28**, 709-723 (2017).
- 1732 20. All natural. *Nature Chemical Biology* **3**, 351 (2007).
- 1733 21. in The "Gold Book", Edn. 2nd. (eds. A.D. McNaught & A. Wilkinson) (Blackwell  
1734 Scientific Publications, Oxford; 1997).
- 1735 22. McLafferty, F.W. Tandem mass spectrometry. *Science* **214**, 280-287 (1981).
- 1736 23. Gross, J.H. in *Mass Spectrometry: A Textbook* 415-478 (Springer Berlin  
1737 Heidelberg, Berlin, Heidelberg; 2011).
- 1738 24. Artyukhin, A.B. et al. Metabolomic "Dark Matter" Dependent on Peroxisomal  $\beta$ -  
1739 Oxidation in *Caenorhabditis elegans*. *Journal of the American Chemical Society*  
1740 **140**, 2841-2852 (2018).
- 1741 25. Edwards, E.D., Woolly, E.F., McLellan, R.M. & Keyzers, R.A. Non-detection of  
1742 honeybee hive contamination following *Vespula* wasp baiting with protein  
1743 containing fipronil. *PLoS One* **13**, e0206385 (2018).
- 1744 26. Hoffmann, T. et al. Correlating chemical diversity with taxonomic distance for  
1745 discovery of natural products in myxobacteria. *Nature Communications* **9**, 803  
1746 (2018).
- 1747 27. Leipoldt, F. et al. Warhead biosynthesis and the origin of structural diversity in  
1748 hydroxamate metalloproteinase inhibitors. *Nat Commun* **8**, 1965 (2017).
- 1749 28. Kang, K.B., Gao, M., Kim, G.J., Choi, H. & Sung, S.H. Rhamnelloides A and B,  
1750 omega-Phenylpentaene Fatty Acid Amide Diglycosides from the Fruits of  
1751 *Rhamnella franguloides*. *Molecules* **23** (2018).
- 1752 29. Remy, S. et al. Structurally Diverse Diterpenoids from *Sandwithia guyanensis*.  
1753 *Journal of Natural Products* **81**, 901-912 (2018).
- 1754 30. Riewe, D., Wiebach, J. & Altmann, T. Structure Annotation and Quantification of  
1755 Wheat Seed Oxidized Lipids by High-Resolution LC-MS/MS. *Plant Physiol* **175**,  
1756 600-618 (2017).
- 1757 31. Senges, C.H.R. et al. The secreted metabolome of *Streptomyces*  
1758 *chartreusis* and implications for bacterial chemistry. *Proceedings of the*  
1759 *National Academy of Sciences* **115**, 2490-2495 (2018).

- 1760 32. van der Hooft, J.J.J. et al. Unsupervised Discovery and Comparison of Structural  
1761 Families Across Multiple Samples in Untargeted Metabolomics. *Anal Chem* **89**,  
1762 7569-7577 (2017).
- 1763 33. Wolff, H. & Bode, H.B. The benzodiazepine-like natural product tilivalline is  
1764 produced by the entomopathogenic bacterium *Xenorhabdus eapokensis*. *PLoS*  
1765 *One* **13**, e0194297 (2018).
- 1766 34. von Eckardstein, L. et al. Total Synthesis and Biological Assessment of Novel  
1767 Albicidins Discovered by Mass Spectrometric Networking. *Chemistry* **23**, 15316-  
1768 15321 (2017).
- 1769 35. Vizcaino, M.I. & Crawford, J.M. The colibactin warhead crosslinks DNA. *Nat Chem*  
1770 **7**, 411-417 (2015).
- 1771 36. Saleh, H. et al. Deuterium-Labeled Precursor Feeding Reveals a New pABA-  
1772 Containing Meroterpenoid from the Mango Pathogen *Xanthomonas citri* pv.  
1773 *mangiferaeindicae*. *J Nat Prod* **79**, 1532-1537 (2016).
- 1774 37. Shannon, P. et al. Cytoscape: a software environment for integrated models of  
1775 biomolecular interaction networks. *Genome research* **13**, 2498-2504 (2003).
- 1776 38. Petras, D. et al. Mass Spectrometry-Based Visualization of Molecules Associated  
1777 with Human Habitats. *Anal Chem* **88**, 10775-10784 (2016).
- 1778 39. Kapon, C.A. et al. Creating a 3D microbial and chemical snapshot of a human  
1779 habitat. *Sci Rep* **8**, 3669 (2018).
- 1780 40. Adams, R.I. et al. Microbes and associated soluble and volatile chemicals on  
1781 periodically wet household surfaces. *Microbiome* **5**, 128 (2017).
- 1782 41. Petras, D. et al. High-Resolution Liquid Chromatography Tandem Mass  
1783 Spectrometry Enables Large Scale Molecular Characterization of Dissolved  
1784 Organic Matter. *Frontiers in Marine Science* **4** (2017).
- 1785 42. Trautman, E.P. & Crawford, J.M. Linking Biosynthetic Gene Clusters to their  
1786 Metabolites via Pathway- Targeted Molecular Networking. *Curr Top Med Chem* **16**,  
1787 1705-1716 (2016).
- 1788 43. Luzzatto-Knaan, T., Melnik, A.V. & Dorrestein, P.C. Mass Spectrometry Uncovers  
1789 the Role of Surfactin as an Interspecies Recruitment Factor. *ACS Chemical Biology*  
1790 (2019).
- 1791 44. Machushynets, N.V., Wu, C., Elsayed, S.S., Hankemeier, T. & van Wezel, G.P.  
1792 Discovery of novel glycerolated quinazolinones from *Streptomyces* sp. MBT27. *J*  
1793 *Ind Microbiol Biotechnol* (2019).
- 1794 45. Yao, L. et al. Discovery of novel xylosides in co-culture of basidiomycetes *Trametes*  
1795 *versicolor* and *Ganoderma applanatum* by integrated metabolomics and  
1796 bioinformatics. *Sci Rep* **6**, 33237 (2016).
- 1797 46. Tripathi, A. et al. Intermittent Hypoxia and Hypercapnia, a Hallmark of Obstructive  
1798 Sleep Apnea, Alters the Gut Microbiome and Metabolome. *mSystems* **3** (2018).
- 1799 47. Smits, S.A. et al. Seasonal cycling in the gut microbiome of the Hadza hunter-  
1800 gatherers of Tanzania. *Science* **357**, 802-806 (2017).
- 1801 48. McDonald, D. et al. American Gut: an Open Platform for Citizen Science  
1802 Microbiome Research. *mSystems* **3**, e00031-00018 (2018).
- 1803 49. Edlund, A. et al. Metabolic Fingerprints from the Human Oral Microbiome Reveal  
1804 a Vast Knowledge Gap of Secreted Small Peptidic Molecules. *mSystems* **2**,  
1805 e00058-00017 (2017).
- 1806 50. Garg, N. et al. Three-Dimensional Microbiome and Metabolome Cartography of a  
1807 Diseased Human Lung. *Cell Host Microbe* **22**, 705-716 e704 (2017).
- 1808 51. McCall, L.I. et al. Mass Spectrometry-Based Chemical Cartography of a Cardiac  
1809 Parasitic Infection. *Anal Chem* **89**, 10414-10421 (2017).
- 1810 52. Watrous, J.D. et al. Directed Non-targeted Mass Spectrometry and Chemical  
1811 Networking for Discovery of Eicosanoids and Related Oxylipins. *Cell Chemical*  
1812 *Biology* (2019).

- 1813 53. Allard, S., Allard, P.M., Morel, I. & Gicquel, T. Application of a molecular networking  
1814 approach for clinical and forensic toxicology exemplified in three cases involving 3-  
1815 MeO-PCP, doxylamine, and chlormequat. *Drug Test Anal* (2018).
- 1816 54. Ernst, M. et al. Did a plant-herbivore arms race drive chemical diversity in  
1817 Euphorbia? *bioRxiv*, 323014 (2018).
- 1818 55. Philippus, A.C. et al. Molecular networking prospection and characterization of  
1819 terpenoids and C15-acetogenins in Brazilian seaweed extracts. *RSC Advances* **8**,  
1820 29654-29661 (2018).
- 1821 56. Li, F., Janussen, D., Peifer, C., Perez-Victoria, I. & Tasdemir, D. Targeted Isolation  
1822 of Tsitsikammamines from the Antarctic Deep-Sea Sponge *Latrunculia biformis* by  
1823 Molecular Networking and Anticancer Activity. *Mar Drugs* **16** (2018).
- 1824 57. Hartmann, A.C. et al. Meta-mass shift chemical profiling of metabolomes from coral  
1825 reefs. *Proc Natl Acad Sci U S A* **114**, 11685-11690 (2017).
- 1826 58. Tobias, N.J. et al. Natural product diversity associated with the nematode  
1827 symbionts *Photobacterium* and *Xenorhabdus*. *Nature Microbiology* **2**, 1676-1685  
1828 (2017).
- 1829 59. Nothias, L.F. et al. Bioactivity-Based Molecular Networking for the Discovery of  
1830 Drug Leads in Natural Product Bioassay-Guided Fractionation. *J Nat Prod* **81**, 758-  
1831 767 (2018).
- 1832 60. Zou, Y. et al. Computationally Assisted Discovery and Assignment of a Highly  
1833 Strained and PANC-1 Selective Alkaloid from Alaska's Deep Ocean. *Journal of the*  
1834 *American Chemical Society* (2019).
- 1835 61. Parkinson, E.I. et al. Discovery of the Tyrobetaine Natural Products and Their  
1836 Biosynthetic Gene Cluster via Metabologenomics. *ACS Chemical Biology* **13**,  
1837 1029-1037 (2018).
- 1838 62. Naman, C.B. et al. Integrating Molecular Networking and Biological Assays To  
1839 Target the Isolation of a Cytotoxic Cyclic Octapeptide, Samoamide A, from an  
1840 American Samoan Marine Cyanobacterium. *Journal of Natural Products* **80**, 625-  
1841 633 (2017).
- 1842 63. Bouslimani, A. et al. Lifestyle chemistries from phones for individual profiling. *Proc*  
1843 *Natl Acad Sci U S A* **113**, E7645-E7654 (2016).
- 1844 64. Schymanski, E.L. et al. Critical Assessment of Small Molecule Identification 2016:  
1845 automated methods. *Journal of Cheminformatics* **9**, 22 (2017).
- 1846 65. Quinn, R.A. et al. Niche partitioning of a pathogenic microbiome driven by chemical  
1847 gradients. *Sci Adv* **4**, eaau1908 (2018).
- 1848 66. Aksenov, A.A., da Silva, R., Knight, R., Lopes, N.P. & Dorrestein, P.C. Global  
1849 chemical analysis of biology by mass spectrometry. *Nature Reviews Chemistry* **1**,  
1850 0054 (2017).
- 1851 67. Tsugawa, H. Advances in computational metabolomics and databases deepen the  
1852 understanding of metabolisms. *Current Opinion in Biotechnology* **54**, 10-17 (2018).
- 1853 68. Johnson, S.R. & Lange, B.M. Open-Access Metabolomics Databases for Natural  
1854 Product Research: Present Capabilities and Future Potential. *Frontiers in*  
1855 *Bioengineering and Biotechnology* **3** (2015).
- 1856 69. Haug, K. et al. MetaboLights--an open-access general-purpose repository for  
1857 metabolomics studies and associated meta-data. *Nucleic Acids Res* **41**, D781-786  
1858 (2013).
- 1859 70. Perez-Riverol, Y. et al. Discovering and linking public omics data sets using the  
1860 Omics Discovery Index. *Nat Biotechnol* **35**, 406-409 (2017).
- 1861 71. Stein, S.E. & Scott, D.R. Optimization and testing of mass spectral library search  
1862 algorithms for compound identification. *Journal of the American Society for Mass*  
1863 *Spectrometry* **5**, 859-866 (1994).
- 1864 72. NIST Standard Reference Database 1A v17.
- 1865 73. Guijas, C. et al. METLIN: A Technology Platform for Identifying Knowns and  
1866 Unknowns. *Anal Chem* **90**, 3156-3164 (2018).

- 1867 74. Horai, H. et al. MassBank: a public repository for sharing mass spectral data for  
1868 life sciences. *J Mass Spectrom* **45**, 703-714 (2010).
- 1869 75. Stravs, M.A., Schymanski, E.L., Singer, H.P. & Hollender, J. Automatic  
1870 recalibration and processing of tandem mass spectra using formula annotation. *J*  
1871 *Mass Spectrom* **48**, 89-99 (2013).
- 1872 76. Wang, J., Peake, D.A., Mistrik, R., Huang, Y. & Araujo, G.D.  
1873 ([http://www.unitylabservices.eu/content/dam/tfs/ATG/CMD/CMD%20Documents/](http://www.unitylabservices.eu/content/dam/tfs/ATG/CMD/CMD%20Documents/posters/PN-ASMS13-a-platform-to-identify-endogenous-metabolites-using-a-novel-high-performance-orbitrap-and-the-mzcloud-library-E.pdf)  
1874 [posters/PN-ASMS13-a-platform-to-identify-endogenous-metabolites-using-a-](http://www.unitylabservices.eu/content/dam/tfs/ATG/CMD/CMD%20Documents/posters/PN-ASMS13-a-platform-to-identify-endogenous-metabolites-using-a-novel-high-performance-orbitrap-and-the-mzcloud-library-E.pdf)  
1875 [novel-high-performance-orbitrap-and-the-mzcloud-library-E.pdf](http://www.unitylabservices.eu/content/dam/tfs/ATG/CMD/CMD%20Documents/posters/PN-ASMS13-a-platform-to-identify-endogenous-metabolites-using-a-novel-high-performance-orbitrap-and-the-mzcloud-library-E.pdf); 2013).
- 1876 77. Sheldon, M.T., Mistrik, R. & Croley, T.R. Determination of ion structures in  
1877 structurally related compounds using precursor ion fingerprinting. *J Am Soc Mass*  
1878 *Spectrom* **20**, 370-376 (2009).
- 1879 78. Sawada, Y. et al. RIKEN tandem mass spectral database (ReSpect) for  
1880 phytochemicals: a plant-specific MS/MS-based data resource and database.  
1881 *Phytochemistry* **82**, 38-45 (2012).
- 1882 79. Smith, C.A., Want, E.J., O'Maille, G., Abagyan, R. & Siuzdak, G. XCMS:  
1883 Processing Mass Spectrometry Data for Metabolite Profiling Using Nonlinear Peak  
1884 Alignment, Matching, and Identification. *Analytical Chemistry* **78**, 779-787 (2006).
- 1885 80. Tautenhahn, R., Patti, G.J., Rinehart, D. & Siuzdak, G. XCMS Online: a web-based  
1886 platform to process untargeted metabolomic data. *Anal Chem* **84**, 5035-5039  
1887 (2012).
- 1888 81. Wanichthanarak, K., Fan, S., Grapov, D., Barupal, D.K. & Fiehn, O. Metabox: A  
1889 Toolbox for Metabolomic Data Analysis, Interpretation and Integrative Exploration.  
1890 *PLOS ONE* **12**, e0171046 (2017).
- 1891 82. Mohimani, H. et al. Dereplication of microbial metabolites through database search  
1892 of mass spectra. *Nature Communications* **9**, 4035 (2018).
- 1893 83. Mohimani, H. et al. Dereplication of peptidic natural products through database  
1894 search of mass spectra. *Nat Chem Biol* **13**, 30-37 (2017).
- 1895 84. Gurevich, A. et al. Increased diversity of peptidic natural products revealed by  
1896 modification-tolerant database search of mass spectra. *Nat Microbiol* **3**, 319-327  
1897 (2018).
- 1898 85. da Silva, R.R. et al. Propagating annotations of molecular networks using in silico  
1899 fragmentation. *PLoS computational biology* **14**, e1006089 (2018).
- 1900 86. Mohimani, H. et al. Automated genome mining of ribosomal peptide natural  
1901 products. *ACS Chem Biol* **9**, 1545-1551 (2014).
- 1902 87. Olivon, F. et al. MetGem Software for the Generation of Molecular Networks Based  
1903 on the t-SNE Algorithm. *Anal Chem* (2018).
- 1904 88. Olivon, F., Roussi, F., Litaudon, M. & Touboul, D. Optimized experimental workflow  
1905 for tandem mass spectrometry molecular networking in metabolomics. *Anal*  
1906 *Bioanal Chem* **409**, 5767-5778 (2017).
- 1907 89. Wehrens, R. et al. Improved batch correction in untargeted MS-based  
1908 metabolomics. *Metabolomics* **12**, 88 (2016).
- 1909 90. Koal, T. & Deigner, H.P. Challenges in mass spectrometry based targeted  
1910 metabolomics. *Curr Mol Med* **10**, 216-226 (2010).
- 1911 91. Bylda, C., Thiele, R., Kobold, U. & Volmer, D.A. Recent advances in sample  
1912 preparation techniques to overcome difficulties encountered during quantitative  
1913 analysis of small molecules from biofluids using LC-MS/MS. *Analyst* **139**, 2265-  
1914 2276 (2014).
- 1915 92. Vuckovic, D. Current trends and challenges in sample preparation for global  
1916 metabolomics using liquid chromatography-mass spectrometry. *Anal Bioanal*  
1917 *Chem* **403**, 1523-1548 (2012).
- 1918 93. Dunn, W.B. et al. Procedures for large-scale metabolic profiling of serum and  
1919 plasma using gas chromatography and liquid chromatography coupled to mass  
1920 spectrometry. *Nature Protocols* **6**, 1060 (2011).

- 1921 94. Taylor, P.J. Matrix effects: the Achilles heel of quantitative high-performance liquid  
1922 chromatography-electrospray-tandem mass spectrometry. *Clin Biochem* **38**, 328-  
1923 334 (2005).
- 1924 95. Annesley, T.M. Ion suppression in mass spectrometry. *Clin Chem* **49**, 1041-1044  
1925 (2003).
- 1926 96. Crüsemann, M. et al. Prioritizing Natural Product Diversity in a Collection of 146  
1927 Bacterial Strains Based on Growth and Extraction Protocols. *Journal of Natural*  
1928 *Products* **80**, 588-597 (2017).
- 1929 97. Wandro, S., Carmody, L., Gallagher, T., LiPuma, J.J. & Whiteson, K. Making It  
1930 Last: Storage Time and Temperature Have Differential Impacts on Metabolite  
1931 Profiles of Airway Samples from Cystic Fibrosis Patients. *mSystems* **2** (2017).
- 1932 98. Zhao, J., Evans, C.R., Carmody, L.A. & LiPuma, J.J. Impact of storage conditions  
1933 on metabolite profiles of sputum samples from persons with cystic fibrosis. *J Cyst*  
1934 *Fibros* **14**, 468-473 (2015).
- 1935 99. Hirayama, A. et al. Effects of processing and storage conditions on charged  
1936 metabolomic profiles in blood. *ELECTROPHORESIS* **36**, 2148-2155 (2015).
- 1937 100. Mushtaq, M.Y., Choi, Y.H., Verpoorte, R. & Wilson, E.G. Extraction for  
1938 metabolomics: access to the metabolome. *Phytochem Anal* **25**, 291-306 (2014).
- 1939 101. Bazsó, F.L. et al. Quantitative Comparison of Tandem Mass Spectra Obtained on  
1940 Various Instruments. *J Am Soc Mass Spectrom* **27**, 1357-1365 (2016).
- 1941 102. Bowen, B.P. & Northen, T.R. Dealing with the unknown: metabolomics and  
1942 metabolite atlases. *J Am Soc Mass Spectrom* **21**, 1471-1476 (2010).
- 1943 103. da Silva, R.R., Dorrestein, P.C. & Quinn, R.A. Illuminating the dark matter in  
1944 metabolomics. *Proc Natl Acad Sci U S A* **112**, 12549-12550 (2015).
- 1945 104. Blaženović, I., Kind, T., Ji, J. & Fiehn, O. Software Tools and Approaches for  
1946 Compound Identification of LC-MS/MS Data in Metabolomics. *Metabolites* **8**  
1947 (2018).
- 1948 105. Ruttkies, C., Schymanski, E.L., Wolf, S., Hollender, J. & Neumann, S. MetFrag  
1949 relaunched: incorporating strategies beyond in silico fragmentation. *J Cheminform*  
1950 **8**, 3 (2016).
- 1951 106. Gerlich, M. & Neumann, S. MetFusion: integration of compound identification  
1952 strategies. *J Mass Spectrom* **48**, 291-298 (2013).
- 1953 107. Böcker, S., Letzel, M.C., Liptak, Z. & Pevukhin, A. SIRIUS: decomposing isotope  
1954 patterns for metabolite identification. *Bioinformatics* **25**, 218-224 (2009).
- 1955 108. Dührkop, K. et al. SIRIUS 4: a rapid tool for turning tandem mass spectra into  
1956 metabolite structure information. *Nat Methods* **16**, 299-302 (2019).
- 1957 109. Dührkop, K., Shen, H., Meusel, M., Rousu, J. & Bocker, S. Searching molecular  
1958 structure databases with tandem mass spectra using CSI:FingerID. *Proc Natl Acad*  
1959 *Sci U S A* **112**, 12580-12585 (2015).
- 1960 110. Tsugawa, H. et al. Hydrogen Rearrangement Rules: Computational MS/MS  
1961 Fragmentation and Structure Elucidation Using MS-FINDER Software. *Anal Chem*  
1962 **88**, 7946-7958 (2016).
- 1963 111. Protsyuk, I. et al. 3D molecular cartography using LC-MS facilitated by Optimus  
1964 and 'ili software. *Nat Protoc* **13**, 134-154 (2018).
- 1965 112. Röst, H.L. et al. OpenMS: a flexible open-source software platform for mass  
1966 spectrometry data analysis. *Nat Methods* **13**, 741-748 (2016).
- 1967 113. Pluskal, T., Castillo, S., Villar-Briones, A. & Oresic, M. MZmine 2: modular  
1968 framework for processing, visualizing, and analyzing mass spectrometry-based  
1969 molecular profile data. *BMC Bioinformatics* **11**, 395 (2010).
- 1970 114. Deutsch, E.W. et al. Proteomics Standards Initiative: Fifteen Years of Progress and  
1971 Future Work. *Journal of Proteome Research* **16**, 4288-4298 (2017).
- 1972 115. Brooksbank, C., Cameron, G. & Thornton, J. The European Bioinformatics  
1973 Institute's data resources. *Nucleic Acids Res* **38**, D17-25 (2010).
- 1974 116. Jones, A.R. et al. The mzIdentML data standard for mass spectrometry-based  
1975 proteomics results. *Mol Cell Proteomics* **11**, M111 014381 (2012).

- 1976 117. Griss, J. et al. The mzTab data exchange format: communicating mass-  
1977 spectrometry-based proteomics and metabolomics experimental results to a wider  
1978 audience. *Mol Cell Proteomics* **13**, 2765-2775 (2014).
- 1979 118. Hoffmann, N. et al. mzTab-M: A Data Standard for Sharing Quantitative Results in  
1980 Mass Spectrometry Metabolomics. *Analytical Chemistry* (2019).
- 1981 119. Wang, M. et al. MASST: A Web-based Basic Mass Spectrometry Search Tool for  
1982 Molecules to Search Public Data. *bioRxiv*, 591016 (2019).
- 1983 120. Scheubert, K. et al. Significance estimation for large scale metabolomics  
1984 annotations by spectral matching. *Nat Commun* **8**, 1494 (2017).
- 1985 121. McDonald, D. et al. The Biological Observation Matrix (BIOM) format or: how I  
1986 learned to stop worrying and love the ome-ome. *GigaScience* **1**, 7 (2012).
- 1987 122. Vazquez-Baeza, Y., Pirrung, M., Gonzalez, A. & Knight, R. EMPeror: a tool for  
1988 visualizing high-throughput microbial community data. *GigaScience* **2**, 16 (2013).
- 1989 123. Bolyen, E. et al. QIIME 2: Reproducible, interactive, scalable, and extensible  
1990 microbiome data science. *PeerJ Preprints* (2018).
- 1991 124. McLafferty, F.W. & Tureček, F.e. Interpretation of mass spectra, Edn. 4th.  
1992 (University Science Books, Mill Valley, Calif.; 1993).
- 1993 125. Viant, M.R., Kurland, I.J., Jones, M.R. & Dunn, W.B. How close are we to complete  
1994 annotation of metabolomes? *Curr Opin Chem Biol* **36**, 64-69 (2017).
- 1995 126. Shahaf, N. et al. The WEIZMASS spectral library for high-confidence metabolite  
1996 identification. *Nature Communications* **7**, 12423 (2016).
- 1997 127. Schymanski, E.L. et al. Identifying small molecules via high resolution mass  
1998 spectrometry: communicating confidence. *Environ Sci Technol* **48**, 2097-2098  
1999 (2014).
- 2000 128. Cleary, J.L., Luu, G.T., Pierce, E.C., Dutton, R.J. & Sanchez, L.M. BLANKA: an  
2001 Algorithm for Blank Subtraction in Mass Spectrometry of Complex Biological  
2002 Samples. *Journal of The American Society for Mass Spectrometry* (2019).
- 2003 129. Demarque, D.P., Crotti, A.E.M., Vessecchi, R., Lopes, J.L.C. & Lopes, N.P.  
2004 Fragmentation reactions using electrospray ionization mass spectrometry: an  
2005 important tool for the structural elucidation and characterization of synthetic and  
2006 natural products. *Natural Product Reports* **33**, 432-455 (2016).
- 2007 130. van der Hoof, J.J.J., Wandy, J., Barrett, M.P., Burgess, K.E.V. & Rogers, S. Topic  
2008 modeling for untargeted substructure exploration in metabolomics. *Proceedings of  
2009 the National Academy of Sciences* **113**, 13738-13743 (2016).
- 2010 131. Olivon, F., Grelier, G., Roussi, F., Litaudon, M. & Touboul, D. MZmine 2 Data-  
2011 Preprocessing To Enhance Molecular Networking Reliability. *Analytical Chemistry*  
2012 **89**, 7836-7840 (2017).
- 2013 132. Winnikoff, J.R., Glukhov, E., Watrous, J., Dorrestein, P.C. & Gerwick, W.H.  
2014 Quantitative molecular networking to profile marine cyanobacterial metabolomes.  
2015 *J Antibiot (Tokyo)* **67**, 105-112 (2014).
- 2016 133. Tsugawa, H. et al. MS-DIAL: data-independent MS/MS deconvolution for  
2017 comprehensive metabolome analysis. *Nat Methods* **12**, 523-526 (2015).
- 2018 134. Marfey, P. Determination of D-Amino Acids .2. Use of a Bifunctional Reagent, 1,5-  
2019 Difluoro-2,4-Dinitrobenzene. *Carlsberg Res Commun* **49**, 591-596 (1984).
- 2020