

GigaScience

From chromatogram to analyte to metabolite. How to pick horses for courses from the massive web-resources for mass spectral plant metabolomics --Manuscript Draft--

Manuscript Number:	GIGA-D-17-00039					
Full Title:	From chromatogram to analyte to metabolite. How to pick horses for courses from the massive web-resources for mass spectral plant metabolomics					
Article Type:	Review					
Funding Information:	<table border="1"> <tr> <td>Conselho Nacional de Desenvolvimento Científico e Tecnológico (246605/2012-0)</td> <td>Mr. Leonardo Perez de Souza</td> </tr> <tr> <td>Max-Planck-Gesellschaft</td> <td>Not applicable</td> </tr> </table>	Conselho Nacional de Desenvolvimento Científico e Tecnológico (246605/2012-0)	Mr. Leonardo Perez de Souza	Max-Planck-Gesellschaft	Not applicable	
Conselho Nacional de Desenvolvimento Científico e Tecnológico (246605/2012-0)	Mr. Leonardo Perez de Souza					
Max-Planck-Gesellschaft	Not applicable					
Abstract:	<p>The grand challenge currently facing metabolomics is the expansion of the coverage of the metabolome from a minor percentage of the metabolic complement of the cell towards the level of coverage afforded by other post-genomic technologies such as transcriptomics and proteomics. In plants this problem is exacerbated by the sheer diversity of chemicals that constitute the metabolome with the number of metabolites in the plant kingdom generally being considered to be in excess of 200 000. In this review we focus on web-resources that can be exploited in order to improve analyte and ultimately metabolite identification and quantification. There is a wide range of available software that not only aids in this but also in the related area of peak alignment, however, for the uninitiated choosing which program to use is a daunting task. For this reason we provide an overview of the pros and cons of the software as well as comments regarding the level of programming skills required to effectively exploit their basic functions. In addition the torrent of available genome and transcriptome sequences that followed the advent of next-generation sequencing has opened up further valuable resources for metabolite identification. All things considered, we posit that only via a continued communal sharing of information such as that deposited in the databases described within the article are we likely to be able to make significant headway towards improving our coverage of the plant metabolome.</p>					
Corresponding Author:	Alisdair Robert Fernie GERMANY					
Corresponding Author Secondary Information:						
Corresponding Author's Institution:						
Corresponding Author's Secondary Institution:						
First Author:	Leonardo Perez de Souza					
First Author Secondary Information:						
Order of Authors:	<table border="1"> <tr><td>Leonardo Perez de Souza</td></tr> <tr><td>Thomas Naake</td></tr> <tr><td>Takayuki Tohge</td></tr> <tr><td>Alisdair Robert Fernie</td></tr> </table>		Leonardo Perez de Souza	Thomas Naake	Takayuki Tohge	Alisdair Robert Fernie
Leonardo Perez de Souza						
Thomas Naake						
Takayuki Tohge						
Alisdair Robert Fernie						
Order of Authors Secondary Information:						
Opposed Reviewers:						
Additional Information:						
Question	Response					
Are you submitting this manuscript to a	Yes					

special series or article collection?	
Please select an option from the menu: as follow-up to "Are you submitting this manuscript to a special series or article collection?"	Functional Metagenomics
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	No
<p>If not, please give reasons for any omissions below.</p> <p>as follow-up to "Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p> <p>"</p>	Review article
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	No

<p>If not, please give reasons for any omissions below.</p> <p>as follow-up to "Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p> <p>"</p>	<p>Review article</p>
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	<p>No</p>
<p>If not, please give reasons for any omissions below.</p> <p>as follow-up to "Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum</p>	<p>Review article</p>

[Standards Reporting Checklist?](#)

"

1 Review:

2
3 **From chromatogram to analyte to metabolite. How to pick horses**
4 **for courses from the massive web-resources for mass spectral plant**
5 **metabolomics**

6 Leonardo Perez de Souza^{1*}, Thomas Naake¹, Takayuki Tohge¹, and Alisdair R. Fernie^{1*}

7 ¹ Max-Planck-Institute of Molecular Plant Physiology, Am Mühlenberg 1, 14476 Potsdam-
8 Golm, Germany

9 * Correspondence: LPerez@mpimp-golm.mpg.de; fernie@mpimp-golm.mpg.de

10 **Abstract**

11 The grand challenge currently facing metabolomics is the expansion of the coverage of the
12 metabolome from a minor percentage of the metabolic complement of the cell towards the
13 level of coverage afforded by other post-genomic technologies such as transcriptomics and
14 proteomics. In plants this problem is exacerbated by the sheer diversity of chemicals that
15 constitute the metabolome with the number of metabolites in the plant kingdom generally
16 being considered to be in excess of 200 000. In this review we focus on web-resources that
17 can be exploited in order to improve analyte and ultimately metabolite identification and
18 quantification. There is a wide range of available software that not only aids in this but also
19 in the related area of peak alignment, however, for the uninitiated choosing which program
20 to use is a daunting task. For this reason we provide an overview of the pros and cons of the
21 software as well as comments regarding the level of programming skills required to effectively
22 exploit their basic functions. In addition the torrent of available genome and transcriptome
23 sequences that followed the advent of next-generation sequencing has opened up further
24 valuable resources for metabolite identification. All things considered, we posit that only via
25 a continued communal sharing of information such as that deposited in the databases
26 described within the article are we likely to be able to make significant headway towards
27 improving our coverage of the plant metabolome.

28 **Keywords:** Arabidopsis, bioinformatics, crop species, GC-MS, LC-MS, peak identification,
29 peak annotation.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

30
31
32
33

34 Background

1
2 35 Metabolomics emerged in the late 1990s with the term coined in a review of Steven Oliver
3 36 [1]. However, the 2000 paper by Fiehn and co-workers wherein gas chromatography (GC)
4 37 coupled to mass spectrometry (MS) defined the chemical composition of a morphological
5 38 and metabolic mutant of the model plant *Arabidopsis thaliana* [2] in doing so they were able
6 39 to describe changes in the level of 326 analytes. This work thus greatly extended on the
7 40 early metabolite profiling study of Sauter et al. [3], which presented the technology as a
8 41 means of putative classification of mode-of-action of pesticides. Thus the advent of
9 42 metabolomics in plants arguably preceded that in microbes and mammals although the
10 43 approach was rapidly adopted in these communities also [2, 4-6]. During the next two
11 44 decades metabolomics had one considerable advantage over profiling technologies such as
12 45 transcriptomics and proteomics in that it is not directly reliant on the genome sequence and
13 46 during this time the species scope of metabolomics rapidly expanded such that it was no
14 47 longer merely a tool for identifying biomarkers of cellular circumstance but additionally one
15 48 of the cornerstones of systems biology and an approach which could provide mechanistic
16 49 insight into metabolic regulation [7-11]. This advantage has subsequently disappeared
17 50 following the widespread adoption of next-generation sequencing and the lack of linear
18 51 relationship between the genome and the metabolome now represents part of the problem
19 52 in identification of unknown analytes [12]. This is nicely exemplified by the fact that
20 53 computation of the size of the metabolome on genome information as attempted by Nobeli
21 54 and co-workers in 2003 for the *E. coli* metabolome and [13] rendered values far smaller
22 55 than the number of metabolites actually measured to date [14]. Whilst the size of the
23 56 metabolome for prokaryotes has been estimated at a couple of thousand, that of the plant
24 57 kingdom dwarves these numbers with estimates ranging between 200 000 and 1 million
25 58 metabolites [15]. Within the last two decades metabolomics has been employed to address
26 59 a wide range of important questions in plant biology including pathway structure [15], the
27 60 influence of metabolism on growth [8, 16], plant ecology [17], various aspects of plant
28 61 genetics including evolution and the domestication syndrome [18-20] as well as detailed
29 62 characterizations of the metabolic response to biotic and abiotic stressors [21, 22].

30
31
32
33
34
35 63 In this review, we discuss two topics. The first is the availability of tools to aid in
36 64 chromatogram evaluation. Since we last reviewed this in 2009 [23], the number of resources
37 65 has exploded as has their diversity in type. In 2009 a number of pathway, analytical
38 66 standards, analytical samples and literature databases were available. In the intervening
39 67 period additional sites providing information on experimental and *in silico* mass
40 68 fragmentation, isotopic labeling, pathway predicted metabolites, integration of
41 69 metabolomics with other platforms and mass spectrometry imaging have become available.
42 70 For each resource we will briefly outline functionality and provide illustrative examples of
43 71 their utility. The second is to review the current status of the broad variety of plant
44 72 metabolomics databases. In this respect we list sources of archived data and their
45 73 respective volumes of data. We also briefly discuss recent meta-analysis which illustrate

74 that despite current hurdles regarding comparability of data there is great potential for
75 cross-study comparisons on metabolite responses in determining common responses
76 between either genetic or environmental perturbations of metabolism. Finally, we will
77 provide an outlook as to how the grand challenge of comprehensivity will best be met and
78 how the power of archived plant metabolic responses best exploited in the future.

79 It is not the scope of this review to discuss the theoretical details of every procedure or to
80 document the subtle differences between the many similar tools referred to here. We
81 rather aim to provide a general idea of the importance and challenges of each step in the
82 metabolomics workflow and to summarize the major functions of each tool while referring
83 to the more comprehensive literature supporting them. We attempt to classify all the
84 resources in a simple and logical manner in order to facilitate understanding of the main
85 functionalities of each one. It is, however, important to mention that while few of the tools
86 presented here provide a complete workflow, most of them are able to perform multiple
87 complementary functions somewhat blurring any initiative to accord their functions specific
88 classifications. Other important information that we include here is how these tools can be
89 accessed. This is usually performed either via command-line-or graphical-user-interface
90 (GUI), the former provides flexibility and facilitating integration, automation and
91 development, while the latter was developed to be intuitive and friendly for unexperienced
92 users. All these features considered allow the researcher to access the information required
93 to choose the “winning horse” under the conditions or “course” in which they are racing.
94 Finally it is also important to highlight that these tools are constantly being update,
95 integrated and discontinued, and while we ensured that all the links provided here were
96 functioning at the time of writing, it is impossible to ensure that to be the case in the future.

97 **Sample preparation and data acquisition**

98 The metabolomics workflow (Figure 1) starts with sample preparation including extraction
99 and often coupled to pre-treatment and chemical derivatization, followed by data
100 acquisition which will depend on the chromatographic system, ionization source and
101 analyzer. Optimization of sample preparation and data acquisition can considerably improve
102 the analysis and is particularly interesting for plant metabolomics where matrix complexity
103 is very high; nevertheless this step is often skipped over in favor of standardization and
104 simplicity which allow for greater sample throughput. Methods for chromatography mass
105 spectrometry based optimization are well developed and usually rely on statistical designs
106 collectively known as Design of Experiments (DoE) [24].

107 While some studies have detailed its application in plant metabolite extraction [25] and
108 liquid chromatography (LC) analysis [26], very few software tools were developed so far
109 focusing on this kind of approach for metabolomics data. That said a couple of interesting
110 software are MUSCLE [27], a tool for the automated optimization of targeted LC-MS/MS
111 analysis that was shown to significantly shorten analysis times and increase analytical
112 sensitivities of targeted metabolite analysis, and FragPred [28], which uses experimental

113 fragmentation from a database to select common fragmentation products that minimize
114 uncertainty about metabolite identities in large-scale MRM experiments, have been
115 published and appear to be highly promising.

116 117 **Data processing**

118 Raw mass spectrometry chromatograms are three dimensional data consisting of a
119 distribution of m/z intensities over the time. Processing this data requires filtering, detecting
120 and integrating relevant features, aligning signals across different samples, extracting
121 compound mass spectra and normalizing the data, all with the final goal of simplifying and
122 hence facilitating data interpretation.

123 Feature detection and peak alignment are the initial steps for extracting information from
124 raw data and corresponds to the process in which relevant signals are identified and
125 quantified across samples, having peak alignment as one of the big challenges to overcome,
126 particularly for LC-MS where retention time is more prone to fluctuations in relation to GC-
127 MS. The many different approaches available to perform these steps of data processing
128 were recently reviewed by [29, 30], and some of the most popular algorithms for feature
129 detection and peak alignment were compared in different works [31, 32]. Most software
130 somehow integrate both steps in the same pipeline to generate a report of signal intensities
131 over samples from raw data, and many of them also include some resource for data analysis
132 and peak annotation that will be discussed later in more detail. In the following section we
133 will detail the available tools for this step, adopting a similar approach in all subsequent
134 sections also (the details of the programs are all given in additional file 1). MetAlign [33] is a
135 versatile tool that performs well with both LC-MS and GC-MS and allows direct conversion
136 from and to vendor formats while most other tools need an extra software for this step. It
137 additionally provides a series of functionalities through other tools that are developed by
138 the same group and integrate directly in the output of MetAlign. XCMS appears to be the
139 most cited software for LC-MS data processing, it was developed for R and implements
140 different algorithms for feature detection and alignment suitable for different kinds of data,
141 while it can be argued that the software requires familiarity with programming and lacks
142 resources for simple data inspection, its platform is, nevertheless, powerful and easily
143 integrated with other tools and its extensive community of users provide a great resource
144 for troubleshooting Moreover, a great number of other tools are built upon the functions of
145 XCMS [34]. Amongst these, TracMass 2 [35], a MATLAB software which provides a GUI in a
146 modular suite, was developed to provide immediate graphical feedback of every step of the
147 processing pipeline, its benchmark paper compared the complexity of different algorithms
148 highlighting the importance of low complexity when dealing with large data files and
149 demonstrating it to be more efficient than MZmine 2 (see below for discussion of this
150 software) and comparable to XCMS, two of the most popular current data processing tools.
151 The particularities of TracMass algorithm makes it more suitable for detecting mass traces in

152 the low mass region that can be missed by other approaches, iMet-Q [36], a C# software
153 with a GUI whose algorithm includes automatic detection of charge state and isotope ratio
154 of detected peaks and was developed to minimize the amount of necessary input
155 parameters significantly facilitates the pipeline for new users. GridMass [37] is a 2D feature
156 detection algorithm implemented in MZmine 2 that is faster than other algorithms and
157 potentially improves detection of low-intensity masses. MSFACTs [38], was one of the first
158 tools developed for peak alignment, it uses peak tables or raw data in the ASCII format as
159 input being limited only to the chromatographic domain, this approach can, however, now
160 be considered outdated when compared with many other resources currently available.
161 MET-IDEA [39] is a more recent and flexible tool, developed by the same group as MSFACTs,
162 for feature detection and alignment with a friendly interface developed in .NET platform. Its
163 features include visualization of integrated peaks and manual integration and display of
164 mass spectra, which can be very helpful for quick data inspection. EasyLCMS [40] is a web
165 application tool with focus on calibration and calculation of targeted metabolite
166 concentration in terms of μmol using algorithms developed for MZmine 2. IDEOM [41] is a
167 metabolomics pipeline using functions from XCMS and MZmatch from an Excel GUI. It also
168 includes automated annotation based on an internal database of exact mass and retention
169 time that can be update by users according to the machine. Massifquant [42] is a feature
170 detection algorithm integrated into XCMS based on a Kalman filter for the detection of
171 isotope trace, this approach was shown to be particularly useful for low-intensity peaks.
172 MET-COFEA [43] is a C++ software accessed via a GUI that implements a novel mass trace
173 based extracted-ion chromatogram extraction that copes better with drifts in the mass
174 trace. It additionally uses compound-associated peak clusters instead of individual features
175 for alignment (this clustering process is an important step to extract metabolite information
176 and simplify data as it will be discussed below). MET-Xalign [44] is an extension for MET-
177 COFEA that can potentially align compounds of samples from different experiments, a hard
178 task for metabolomics datasets that is not approached by most other tools. apLCMS [45], is
179 an R package for high mass accuracy LC-MS, which tries to be user friendly by providing a
180 file-based operation and a wrapper function for a single command line batch process of LC-
181 MS data, however, still requires quite some computational knowledge to operate.
182 xMSanalyzer [46] is an R package for improving feature detection that integrates with
183 existing packages such as apLCMS and XCMS, it systematically re-extracts features with
184 multiple parameter settings and merges data to optimize sensitivity and reliability. Yamss
185 [47] is a recently developed R package focused in providing high-quality differential analysis
186 implementing a method based on bivariate approximate kernel density estimation for peak
187 identification. In addition to the tools mentioned above there are a few tools for data
188 processing that exclusively perform peak detection or alignment such as peak-grouping-
189 alignment [48], an approach where information from grouping peaks within samples
190 improve alignment across samples, and PTW [49] a fast alignment algorithm based on a
191 variation of parametric time warping working on detected features rather than on complete
192 profile data. In addition, cosmiq

193 (<http://www.bioconductor.org/packages/devel/bioc/html/cosmiq.html>) is a peak detection
194 algorithm to improve detection of low abundant signals that can be easily integrated with
195 XCMS. These algorithms represent an important effort in improving the existing approaches
196 but they are much less accessible since they need to be integrated with other tools that
197 usually perform similar functions and in some instances this requires quite advanced
198 computational skills.

199 It is important to note the significant differences between GC-MS and LC-MS which are
200 intrinsic to the features of each system, and can be summarized as a much higher efficiency
201 and stability in GC over LC separation followed by a very stable fragmentation in traditional
202 GC ion sources in contrast with the typical atmospheric pressure ionization employed with
203 LC. This significantly influences the processes of peak alignment and spectra annotation, and
204 while most of the tools developed with a focus towards LC-MS can also be used for
205 processing GC-MS data, there are many developed with a particular focus on processing GC-
206 MS data, making use of different strategies for peak alignment and integrating metabolite
207 annotation by matching spectra to libraries. AMDIS [50], developed with the support of U.S.
208 Department of Defense, is one of the most popular GC-MS processing tools, it automatically
209 extracts component mass spectra from GC-MS data and uses it for search in mass spectral
210 libraries, a disadvantage of this software is that the output requires extensive treatment to
211 be used for further analysis. However Metab [51], an R package based on functions of XCMS
212 was developed to automate the pipeline for analysis of GC-MS data processed by AMDIS
213 dealing with the issue of its output data. MetaQuant [52] is a tool that uses retention index
214 to define metabolites but it depends on other deconvolution software like AMDIS to extract
215 mass spectra. Both MetaboliteDetector [53] and TagFinder [54] provide an efficient pipeline
216 performing deconvolution, peak detection, compound identification, alignment based on
217 Kovats retention index using alkane mix and quantification, and provide an interactive user
218 interface facilitating use by unexperienced users. They do however require several manually
219 input and data check steps that are time consuming and negate truly high throughput.
220 TargetSearch [55] uses similar approaches to process data, identify and quantify targeted
221 metabolites based on retention time index and spectra matching of multiple correlated
222 masses but it is highly automated and efficient thus allowing the processing of large sample
223 sets. PyMS [56] is an alternative to the previously mentioned interactive software, providing
224 similar functions but being particularly suitable for scripting of customized processing
225 pipelines and for data processing in batch mode working in Python. MET-COFEI
226 (<http://bioinfo.noble.org/manuscript-support/met-cofei/>) uses reconstructed compound
227 spectra instead of individual peaks to align signals across samples, which is expected to
228 improve peak information for downstream analyses, it also match spectrum against an user-
229 specific library. TNO-DECO [57] uses a segmentation approach to allow the performance of
230 simultaneous deconvolution of multiple chromatographic MS files in a semi-automated
231 fashion in MATLAB, thereby eliminating peak alignment. By contrast, MetaMS [58] is a
232 pipeline for high-throughput GC-MS processing based on XCMS for peak detection and

233 alignment and CAMERA for compound spectra extraction which is annotated based on
1 234 match with a database, this tool may be convenient for users that already implement XCMS
2 235 analysis of other data, but this kind of processing is not optimal for GC-MS when compared
3 236 with other processing types. Maui-VIA [59] implements a graphical interface that facilitates
4 237 visual inspection of identifications and alignments providing faster interaction with the data.
5 238 eRah [60] is an R tool that integrates a novel spectral deconvolution method using
6 239 multivariate techniques based on blind source separation, alignment of spectra across
7 240 samples without the need of internal standards for calculating retention indexes,
8 241 quantification, and automated identification of metabolites by spectral library matching, in a
9 242 fully automated pipeline, even though internal standards are not necessary they are still
10 243 recommended to increase reliability in metabolite identification. The software ADAP-GC 3.0
11 244 [61] uses a deconvolution algorithm based on hierarchical clustering of fragment ions, the
12 245 updated version is incorporated into the MZmine 2 platform and addressed issues from the
13 246 first version such as fragment ions that are produced by more than one co-eluting
14 247 components, and improved sensitivity and robustness. Finally, MetPP [62] is a processing
15 248 tool that includes normalization and statistical analysis but is directed towards data
16 249 emanating from GC×GC-TOF MS system.

25
26 250 Extracting compound mass spectra is another important step of data processing that
27 251 reduces data complexity by many orders of magnitude by identifying m/z signals that belong
28 252 to the same compound and provide essential information for further metabolite annotation
29 253 through the reconstructing of mass spectra. While this process is usually integrated in GC-
30 254 MS tools for feature detection, alignment and annotation, as mentioned above, there are
31 255 many approaches to deal with LC-MS data such as the ones employed by CAMERA [63] a
32 256 package developed in R to extract compound spectra, annotate isotopes and adducts, and
33 257 propose compound mass as an extension to XCMS, it is easy to use in combination with this
34 258 software and provides a significant reduction on data complexity. AStream [64] is another R
35 259 package very similar to CAMERA but using a simpler algorithm for grouping the peaks.
36 260 ALlocator [65], is a web based workflow that applies centwave from XCMS for feature
37 261 detection followed by spectra deconvolution either by CAMERA or by the ALlocatorSD
38 262 algorithm which is optimized for dealing with the particularities of ^{13}C labeled data by
39 263 grouping mirrored isotopes (lighter isotopologues from feeding experiment). MSClust [66],
40 264 has the same general features as the others but it was developed in the C++ language and it
41 265 is optimized to work with the output files of MetAlign. RAMClustR [67] was developed in
42 266 MATLAB and implemented in R, accepting directly the output of XCMS. The authors suggest
43 267 the use of a workflow consisting of data acquisition under both low and high collision energy
44 268 as a way to improve the quality of the spectra generated by feature clustering and provide a
45 269 data format that can be submitted directly to the MassBank Database and NIST MSSearch
46 270 program. By contrast, RAMSY [68] uses average peak ratios and their standard deviations
47 271 rather than correlation to allow the recovery of compound spectra, the performance of this
48 272 approach is typically better than the results from correlation methods, furthermore, the
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

273 script for MATLAB is available or it can be run from a web interface with a .csv table as
274 input.

275 The last step of data processing, data normalization, is essential for further data analysis in
276 order to remove bias introduced by sample preparation from meaningful biological
277 variation. Most methodologies rely either on the use of internal standards statistical means
278 for normalization. Most data normalization procedures are usually integrated in data
279 analysis tools, but there are few examples of more specialized tools such as MetTailor [69]
280 that uses a dynamic block summarization method for correcting misalignments reducing
281 missing data and apply an RT-based local normalization procedure, or Normalyzer [70] that
282 uses twelve different well known normalization methods and compares the results based on
283 different parameters. IntCor [71] that corrects for peak intensity drift effects based on
284 variance analysis, MetNormalizer [72] which allows normalization and integration of
285 multiple batches in large scale experiments using support vector regression, and EigenMS
286 [73] which detect bias trends in the data and eliminates them using single value
287 decomposition are also highly useful. All of these software are implemented in R, however,
288 with the exception of Normalyzer which can be also used in a web interface they all require
289 considerable familiarity with this programming language.

290 A common feature of mass spectrometry data is the presence of multiple peaks for
291 individual fragments resulting from the distribution of natural isotopes which are
292 particularly interesting and explored in stable isotope labeling experiments. There are a few
293 tools for correcting and extracting label enrichment from processed data such as Corrector
294 [74], IsoCor [75] and ICT [76]. These tools are very similar all being based on the same
295 matrix calculation. Corrector was developed to work on the output of TagFinder but data
296 processed with most other tools can be easily arranged in a similar table format. IsoCor
297 provides a GUI with a few different options including corrections for the label input whereas
298 ICT includes features to process data from tandem MS. Nevertheless most data processing
299 pipelines available are not particularly efficient for dealing with this kind of experiment, to
300 fill this gap there are some specialized tools like mzMatch-ISO [77], integrated in the
301 mzMatch pipeline. This software is capable of targeted and untargeted processing of
302 labeled datasets and the output includes a set of plots summarizing the pattern of labelling
303 observed per peak allowing users to quickly explore data. MetExtract [78] which relies on a
304 mixture of cultures from the same organism under natural and labeled media to select
305 signals that show a clear pattern of isotopic enrichment. However, the approach requires
306 the labeled fraction to be fully labeled and the tracer to be highly pure to get the proper
307 isotopic distributions. X13CMS [79] and geoRge [80], both run on the R platform using GC-
308 MS output, the former algorithm iterating over MS signals in each mass spectra using the
309 mass difference due to the label, while the latter uses statistical testing to distinguish
310 Spectral peaks originated from labeled metabolites resulting in significant less false
311 positives. The MIA program [81] detects isotopic enrichment in GC-MS datasets in a non-
312 targeted manner, providing an easy GUI to visualize mass isotopomer distributions (MID) of

1 313 the detected fragments as barplots including confidence intervals and quality measures,
2 314 tools for differential analysis of relative mass isotopomer abundance across samples and
3 315 network assembly based on pairwise similarity of MID that can reveal related metabolites.

4
5 316 Another important feature of many mass spectrometry systems is their capability of
6
7 317 performing tandem mass spectrometry. While this can significantly improve data in many
8 318 ways, it adds another level of complexity for data processing. A very common use of tandem
9
10 319 MS is to increase selectivity and accuracy in targeted analysis and MRManalyzer [82],
11 320 MMSAT [83] and MRMPROBS [84] are useful tools developed for processing data from
12 321 multiple reaction monitoring experiments. MMSAT [83] is a web tool that takes mzXML files
13 322 as the input, it is able to automatically quantify MRM peaks but lacks metabolite
14 323 identification capability. By contrast, MRMPROBS [84] detects and identifies metabolites
15 324 automatically, providing a user-friendly GUI for data analysis. The algorithm has one
16 325 limitation that it needs at least two transitions per metabolite in order to discriminate the
17 326 target metabolite from isomeric metabolites and the background noise. Similarly,
18 327 MRManalyzer [82] is an R tool allowing processing, alignment, metabolite identification,
19 328 quality control check and statistical analysis of large datasets that transforms data in
20 329 “pseudo” accurate m/z, in order to use the centwave algorithm from XCMS for peak
21 330 detection. Untargeted metabolomics analysis can also take advantage of tandem MS,
22 331 particularly for compound annotation, and there are few resources for dealing with the
23 332 complexity of such experiments such as decoMS2 [85], an R package for deconvoluting MS2
24 333 spectra eliminating contaminating fragments without the need of sacrificing sensitivity in
25 334 favor of sensibility by narrowing the window of isolation for collision-induced dissociation
26 335 (CID) during data acquisition. This approach requires MS2 data to be acquired under low
27 336 and high collision energies to solve the mathematical equations potentially reducing
28 337 sensitivity of the method. Similarly MS2Analyzer [86] is a java software for identifying
29 338 neutral losses, precursor ions, product ions and m/z differences from MS2 spectra based on
30 339 a list of predefined transitions. These features are essential for structure elucidation using
31 340 mass spectrometry and the software provides a fast and high-throughput platform for
32 341 extracting this data. MS2LDA [87] is based on latent Dirichlet allocation (LDA), an algorithm
33 342 originally used for text mining that was adapted to generate a list with blocks of co-
34 343 occurring fragments and losses providing results similar to MS2Analyzer but without the
35 344 need of user specified precursor/product transitions. MS-DIAL [88] and MetDIA [89] both
36 345 deal with Data-independent acquisition (DIA) data, an interesting approach for untargeted
37 346 metabolomics that acquire MS2 spectra for all precursor ions simultaneously with the
38 347 complication that it uses larger isolation windows, hence increasing the probability of
39 348 contamination in the MS2, and it loses the relation between precursor and fragment ions.
40 349 MS-DIAL addresses these problems by a mathematical deconvolution based on GC-MS
41 350 processing tools in a fully untargeted manner, whilst achieving the metabolite identification
42 351 through a spectrum-centric library matching. By contrast, MetDIA [89] uses algorithms from
43 352 XCMS for peak detection and alignment combined with a targeted approach based on
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

353 matching metabolites in a library to the detected peaks, thus achieving higher sensitivity
354 and specificity on metabolite identification and wider metabolite coverage.

355 A trade-off for most of the more flexible and powerful resources presented here is that they
356 have multiple parameters that need to be optimized, and recently a number of tools try to
357 assist in evaluating and automatizing this process. In this context IPO [90] was developed to
358 perform automatic optimization of XCMS parameters based on design of experiment ,
359 Credentialing Features [91] optimize detection based on regular and ¹³C-enriched ,
360 MetaboQC [92] is a quality control approach that evaluates alignment and suggests optimal
361 parameters for feature detection based on discrepancies between replicate samples , and
362 SIMAT [93] allows the selection of the optimal set of fragments and retention time windows
363 for target analytes in GC-SIM-MS based analysis.

364 **Data analysis**

365 Metabolomics datasets are usually characterized by high dimensionality, heteroscedasticity
366 (i.e. the variance in errors is not constant across the dataset) and differences of orders of
367 magnitude across metabolite concentrations and fold changes, making it challenging to
368 extract and visualize useful information from processed data. There are numerous
369 approaches for data scaling, reduction, visualization and statistical analysis particularly
370 useful for analyzing metabolomics data, many of them very well established such as analysis
371 of variance (ANOVA), hierarchical cluster analysis (HCS), principal component analysis (PCA)
372 and partial least squares discriminant analysis (PLS-DA) to mention just a few. There are
373 many general statistical software capable of performing most of these functions, but also a
374 variety of software tools exist combining procedures relevant to metabolomics in a single
375 pipeline and thus facilitating the workflow such as DeviumWeb
376 (<https://github.com/dgrapov/DeviumWeb>), BioStatFlow (<http://biostatflow.org/>),
377 MetaboLyzer [94], metaP-Server [95], Fusion ([https://fusion.cebitec.uni-](https://fusion.cebitec.uni-bielefeld.de/Fusion/login)
378 [bielefeld.de/Fusion/login](https://fusion.cebitec.uni-bielefeld.de/Fusion/login)) , Pathomx [96], MSPrep [97], MixOmics (<http://mixomics.org/>)
379 and COVAIN [98].

380 Other interesting and somehow more specialized tools include RepExplore [99] which
381 exploits information from technical replicate variance to improve statistics of differential
382 expression and abundance of omics datasets, KMMDA [100] and Metabomxtr [101] which
383 deal with the troublesome issue of missing metabolite values, the former through a kernel-
384 based score test and the later through mixed-model analysis. Similarly, PeakANOVA [102]
385 identifies peaks that are likely to be associated with one compound and uses them to
386 improve accuracy of quantification, a particularly useful approach for experiments with
387 limited sample size. SPICA [103], is a tool that aims at extracting relevant information from
388 noisy data sets by analyzing ion-pairs instead of individual ions. MetabR [104], normalizes
389 data using linear mixed models and tests for treatment effects with ANOVA. By contrast
390 MPA-RF [105] combines random forests with model population analysis for selecting
391 informative metabolites. Qcscreen [106], helps to verify data consistency, measurement

392 precision and stability of large scale biological experiments. The program SpectConnect
393 [107] identifies conserved metabolites in GC-MS datasets. Finally, MathDAMP [108], a
394 Mathematica package for Differential Analysis of Metabolite Profiles highlights differences
395 within raw LCMS and GCMS dataset.

396 **Metabolite annotation**

397 Metabolite annotation is often considered the most challenging step and as such represents
398 a major bottleneck for metabolomics studies. Even though the gold standard for structural
399 characterization remains NMR characterization of the pure compound [109, 110], MS based
400 metabolomics offers many advantages including lower cost, higher sensitive and
401 throughput, and it can be easily hyphenated with chromatography while still providing
402 considerable structural information. As a consequence great efforts have been made to
403 improve mass spectrometry based metabolite annotation, and a battery of interesting tools
404 were developed with this goal in mind. Structural information is normally extracted from
405 mass of molecular ion in high-resolution MS (HRMS) which can provide the molecular
406 formula and fragmentation pattern. It is important to note that most strategies for
407 metabolite annotation rely heavily on information retrieved from databases of molecular
408 formulas, spectra and pathways which will be discussed in more detail below.

409 The most common tools are based on matching spectra or exact masses from unknown
410 compounds against spectral data deposited in some database. One example using this
411 approach is MetaboSearch [111], which accepts either a list of m/z or the output of CAMERA
412 as input and searches against four major metabolite databases, Human Metabolome
413 DataBase (HMDB), Madison Metabolomics Consortium Database (MMCD), Metlin, and
414 LipidMaps. Similarly, PUTMEDID-LCMS [112] developed in the Taverna Workflow
415 Management System, also integrates a step of compound mass spectra extraction to define
416 a molecular formula from high resolution m/z that is then matched against a predefined list
417 of molecular formulas to annotate compounds. MetAssign [113] is integrated in mzMatch
418 and it considers the uncertainty related with metabolite annotation using a Bayesian
419 clustering approach to assign peak groups, this approach has the advantage of providing a
420 quantitative values for uncertainty/confidence in the outputs that can be used in further
421 analysis. The program SIRIUS [114] is a Java-based software that combines high accuracy
422 mass with isotopic pattern analysis to distinguish even molecular formulas in higher mass
423 regions. Furthermore it also analyses the fragmentation pattern of a compound using
424 fragmentation trees that can be directly uploaded to CSI:FingerID (described below) via a
425 web service. MFSearcher [115] is a tool that efficiently searches high accuracy masses
426 against a database of pre-calculated molecular formulas with fixed kinds and numbers of
427 atoms that are further queried against different databases, HR3 [116] is a similar tool for
428 molecular formula calculation and query in external databases. It uses different sets of rules
429 for heuristic filtering of candidate formulas instead of a pre-calculated database which
430 makes it slightly slower than MFSearcher, but HR3 includes compounds with atoms that are

431 not present in MFSeacher's list as well as considering matches to the isotopic pattern within
1 432 its annotations. Another level of biologically relevant information is added by many tools
2 433 that incorporate pathway information to assist annotation and interpretation of results such
3 434 as Metabolome searcher [117], a web-based application to directly search genome-
4 435 constructed metabolic databases which includes MetaCyc with data on plant metabolism.
5 436 MassTRIX [118] is a web interface that takes a mass peak list from HRMS as input and
6 437 matches them against KEGG compounds database returning a pathway map with the
7 438 matches, organisms can be selected and the output represents organism-specific and extra-
8 439 organism items differentially colored to assist interpretation. MetabNet [119] is an R
9 440 package to perform targeted metabolome wide association study of specific metabolites,
10 441 this approach uses the correlation of all mass signals with the targeted metabolite across
11 442 samples to build networks that can be visualized in pdf or exported to Cytoscape. This can
12 443 be a very useful approach to identify related compounds and associate them to metabolic
13 444 pathways. Similarly, ProbMetab [120] is an R package for probabilistic annotation of
14 445 compounds based on the method developed by Rogers et al. (2009) [121] that incorporates
15 446 information on possible biochemical reactions between the candidate structures to assign
16 447 higher probabilities to compounds that form substrate/product pairs within the same
17 448 sample. MI-Pack [122], implemented in python, calculates differences in mass between all
18 449 molecular formulas annotated from HRMS and compares them to known substrate/product
19 450 pairs from KEGG, but matches are considered based on the error between experimental and
20 451 theoretical masses compared to a threshold defined by a calculated mass error surface.
21 452 PlantMAT [123] is a particularly interesting tool specifically for the investigation of plant
22 453 specialized metabolism, which uses an approach based on common metabolic building
23 454 blocks to predict combinatorial possibilities of phytochemical structures used for annotation
24 455 and as such is a highly effective way to search the chemical space surrounding a (set of)
25 456 metabolite(s)

38
39 457 Another more recent and promising approach made possible by the huge amount of data
40 458 available uses algorithms, mostly based on machine learning, to predict molecular
41 459 properties of unknown compounds from its tandem mass spectra. All the tools listed below
42 460 provide similar web interfaces for putative metabolite identification differing mainly on the
43 461 algorithms used to perform the identification and the overall performance. MetFrag [124]
44 462 retrieves candidate structures either from databases based on exact mass or from user
45 463 specified SDF files, fragments them using a bond dissociation approach and compares the
46 464 fragments with the input spectra scoring matches based on a series of rules. The candidates
47 465 can also be filtered to facilitate the analysis based on relevant factors such as metabolite
48 466 origin, composition, LC retention time and metadata from the databases. Besides the Java
49 467 web-interface a command line version and an R package are provided which are more
50 468 suitable for batch processing and integration with other tools. In a very similar approach
51 469 MolFind [125] retrieves candidates from databases based on exact mass, filters them by
52 470 comparing experimentally measured retention index, ECOM50 (the energy in eV required to
53
54
55
56
57
58
59
60
61
62
63
64
65

471 fragment 50% of a selected precursor ion) and drift time (for ion mobility MS) with
1 472 predicted ones, and analysis CID of the best candidates using MetFrag. CFM-ID [126] is
2 473 based on competitive fragmentation modeling, a probabilistic generative model that uses
3 474 machine learning to learn its parameters from data. It can be used to predict spectra of
4 475 known chemical structures, to annotate peaks in the spectra of a known compound or to
5 476 predict candidate structures for an unknown compound by ranking candidates in terms of
6 477 how closely the predicted spectra match the input. MAGMa [127], extends prediction based
7 478 on substructure assignment by creating hierarchical trees of predicted substructures
8 479 capable of explaining MSⁿ data, where each level takes into account the restrictions
9 480 imposed by the assignment of precursor and subsequent fragmentation. FingerID [128]
10 481 developed a model based on a large dataset of tandem MS from MassBank and uses a
11 482 support vector machine to predict the molecular fingerprint of the unknown spectra and
12 483 compare this with the fingerprint of compounds in a large molecular database. CSI:FingerID
13 484 [129] is a more recent tool based on fingerID that includes computation of fragmentation
14 485 tree achieving the best search performance so far. Besides the web interface it can be also
15 486 queried directly through Sirius but it currently does not support batch mode. Finally
16 487 MetFusion [130] is a Java web tool that combines spectra database matching against
17 488 MassBank with the prediction based annotation provided by MetFrag.

27 489 **Data interpretation**

28 490 Interpretation of omics data is usually complicated by the amount and complexity of data.
29 491 There are many tools to assist metabolomics data interpretation, particularly for its
30 492 visualization by mapping metabolites into pathways and providing biological context, and
31 493 for the integration with data from different platforms (e.g. transcriptomics, proteomics see
32 494 Tohge et al. (2015) [15] for details). As for metabolite annotation, these tools usually rely
33 495 upon knowledge stored in metabolite and pathway databases, and many of them include
34 496 some kind of statistical analysis such as pathway enrichment and correlation analysis.

35 497 Visualization tools provide a simple mean of representing and mapping metabolic changes
36 498 in tools like PATHOS [131], PathWhiz [132] and iPath [133]. They can often provide some
37 499 kind of pathway structure analysis such as PathVisio [134], FunRich [135], BiNChE [136] and
38 500 MPEA [137] that uses pathway enrichment analysis and PAPI [138] that calculates pathway
39 501 activity scores to represent the potential metabolic pathway activities, and performs
40 502 statistical analysis to investigate differences in activity between conditions. Tools like
41 503 InCroMAP [139], IIS [140], KaPPA-View4 [141], MapMan [142], ProMeTra [143] which is
42 504 integrated with MeltDB 2.0, Paintomics [144], VANTED [145], MBROLE [146] and IMPaLA
43 505 [147] go one step further and integrate metabolomics processed data with other omics
44 506 platforms, particularly transcriptomics, providing analysis and visualization of large
45 507 integrated datasets to assist data interpretation.

46 508 Few tools try to actually use mass spectra features to build the networks, which can also
47 509 improve annotation of unknown compounds. MetaNetter [148] uses raw high-resolution
48 510

510 data and a list of potential biochemical transformations to infer metabolic networks.
1 511 MetaMapR [149] builds chemical and spectral similarity networks based on annotated and
2 512 unknown compounds. ChemTreeMap [150] uses annotated structures and a computational
3 513 approach to produce hierarchical trees based on compound similarity to assist visualization
4 514 of chemical overlap between molecular datasets and the extraction of structure–activity
5 515 relationships. MetFamily [151], groups metabolites in families based on an integrated
6 516 analysis of MS1 abundances and MS/MS facilitating further data interpretation. MetCirc
7 517 (<https://www.bioconductor.org/packages/release/bioc/html/MetCirc.html>) is an R tool
8 518 particularly useful for comparative analysis from cross-species and cross-tissue experiments
9 519 through computation of similarity between individual MS/MS spectra and visualization of
10 520 similarity based on interactive graphical tools, and TrackSM [152] is a Java tool that uses
11 521 molecular structure similarities to assign newly identified biochemical compounds to known
12 522 metabolic pathways.

20 523 **Databases**

21
22 524 It must be clear from previous sections that mass spectrometry based metabolomics,
23 525 particularly metabolite annotation and data interpretation, relies heavily upon data from
24 526 characterized mass spectra, molecular properties of analytes and metabolic pathways.
25 527 While all the different techniques offer a lot of flexibility, metabolomics struggles with
26 528 standardization and a great volume of metadata when compared with other omics
27 529 techniques and still lags behind most of them in terms of public repositories of published
28 530 data. Nonetheless there are a wealth of databases with useful information for mass
29 531 spectrometry based plant metabolomics and we try to summarize some of the most
30 532 relevant and the structure and functionalities of resources available.

31
32
33 533 Chempidder [153], PubChem [154], ChEBI [155], ChEMBL [156], ChemBank [157], HMDB
34 534 [158], MMCD [159] and MMsINC [160] are all large databases of small molecules with
35 535 information such as chemical structure, molecular formula and molecular/exact mass, many
36 536 of these databases complement each other and data exchange between them is very
37 537 common, nevertheless it is important to be aware of the sources of data in each one of
38 538 them and to which extent these data is curated, Chempidder for instance has more than 58
39 539 million structures automatically retrieved from over 450 different sources, with only a
40 540 fraction of this being manually curated by registered users while the majority of data only
41 541 went through some sort of automatic curation and elimination of redundant entries.
42 542 Overall such huge databases are particularly useful for looking for physico-chemical
43 543 properties of identified metabolites and checking for possible candidates based solely on
44 544 their mass.

45
46
47 545 There are a few plant specific databases with curated information on chemical composition
48 546 and distribution across different plant species as well, namely KNApSACK [161] with
49 547 information of more than 50,000 metabolites, chemical composition of over 22,000 species,
50 548 Flavonoid viewer [162] with 6,902 molecular structures of flavonoids from 1,687 plant
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

549 species, Dr. Duke's Phytochemical and Ethnobotanical Databases
1 550 (<https://phytochem.nal.usda.gov/phytochem/search>) with information on 29,585 chemicals
2
3 551 of 3,686 medicinal plants, BioPhytMol [163] a resource on anti-mycobacterial
4 552 phytomolecules and plant extracts holding 2,582 entries including 188 plant families,
5
6 553 comprised of 692 genera and 808 species, and 633 active compounds and plant extracts
7 554 identified against 25 target mycobacteria, and EssOilDB [164] with 123,041 essential oil
8
9 555 records from 92 plant families. These are very interesting resources for screening chemical
10 556 composition of specific species and analyzing chemical distribution species wide, and all of
11 557 the data in these databases is manually curated. From all this resources KNApSack is
12
13 558 particularly useful not only for the larger amount of data but also for providing an easy
14
15 559 platform to access and extract information quickly.

16
17 560 Databases providing mass spectra of pure compounds under controlled conditions
18 561 developed to allow search for common spectra features for the identification of unknown
19 562 compounds are an essential resource for MS based identification of metabolites. As
20 563 previously mentioned the great stability and reproducibility of GC-MS generates reliable
21 564 fragmentation patterns and relative retention indexes that are very efficient for metabolite
22 565 annotation by spectra matching. NIST is a very popular commercial library for GC-MS
23 566 annotation, that also provide free access to some data through NIST Chem WebBook
24 567 (<http://webbook.nist.gov/chemistry/>), containing mass spectra of 33,000 compounds. SDDBS
25 568 (http://sdbs.db.aist.go.jp/sdbs/cgi-bin/cre_index.cgi) with 25,000 mass spectra is the
26 569 database from the National Institute of Advanced Industrial Science and Technology (AIST)
27 570 from Japan. Both of them are limited in the fact that they do not offer an interface for
28 571 spectra matching and the user have limited access to data, so those are only useful for
29 572 checking the spectra of targeted compounds. Some more interesting freely-accessible plant
30 573 specific GC-MS libraries include the Golm metabolome database [165] with a total of 26,590
31 574 spectra and 4,663 analytes at the time this article was written and the VocBinBase [166]
32 575 includes 1,537 unique mass spectra at the time this article was written. Both of these
33 576 databases can be downloaded and integrated to processing tools for metabolite annotation
34 577 based on spectra matching.

35 578 One of the greatest efforts in the field of metabolomics has been directed to the
36 579 development of databases of mass spectra obtained from LC-MS analysis. The higher
37 580 flexibility of this technique compared to GC-MS in terms of the chemical space that it can
38 581 analyze comes with the drawback of a high sensitivity to multiple factors that can influence
39 582 mass spectra quality and reproducibility. LC-MS databases are usually characterized by the
40 583 greatest volume of metadata that accompanies the analytical data, and a more complex
41 584 structure for search based in spectra features when compared to GC-MS databases. Some
42 585 large general LC-MS databases include MassBank [167], a public repository of mass spectra
43 586 with 41,092 spectra of 15,828 compounds obtained by 26 different systems (at the time of
44 587 writing). This database is very accessible allowing search by submitted spectra or simply by
45 588 typing in spectral features, mass or targeted compound name, it furthermore allows users

589 to directly extract spectra during data processing through many tools like RAMClustR,
1 590 RMassBank and Mass++. METLIN [168] currently contains 961,829 molecules from which
2 591 200,000 have in silico MS/MS data, and 14,000 mass spectra at multiple collision energies in
3 592 positive and negative ionization mode. METLIN also integrates isoMETLIN [169] that allows
4 593 the search of isotopologues for all METLIN metabolites based on m/z and isotopes of
5 594 interest, and includes experimental data on hundreds of isotopic labeled metabolites that
6 595 can be used to obtain information of precursor atoms in the fragments, both databases can
7 596 be accessed after free registration and searching by mass is fast and easy with the
8 597 advantage that it allows the user to select possible adducts and spectra conditions and
9 598 search directly the mass observed in the spectra. T3DB [170], is a database for toxin data,
10 599 many of which are plant secondary metabolites, with MS, MS-MS and GC-MS spectra of
11 600 3,600 common toxic substances (at the time of writing). mzCloud is a new database with a
12 601 more complex organizing structure that can improve and facilitate data interpretation,
13 602 currently with 6,255 compounds analyzed in different conditions totalizing 1,913,621
14 603 spectra arranged in 9,896 tree structures. It allows the user to easily navigate through
15 604 different spectra of a single compound through its tree structure and also includes
16 605 visualization of predicted molecular formula of the fragments in the spectra
17 606 (<https://www.mzcloud.org/>). Finally the recently developed MoNA
18 607 (<http://mona.fiehnlab.ucdavis.edu/>) is intended to be a centralized, collaborative database
19 608 of metabolite mass spectra and metadata, currently containing over 200,000 mass spectral
20 609 records from experimental and in-silico libraries from different sources. The search is limited
21 610 to name, compound class, molecular formula or exact mass of the metabolite, it can be
22 611 filtered by type of spectra, and the results are presented as a single list of individual
23 612 interactive spectra next to the metadata making it easy to navigate through different
24 613 spectra. The great diversity of phytochemicals observed in plants represent an important
25 614 portion of all these numbers, and a few plant specific databases are available such as
26 615 Spektraris [171], a LC-MS of about 500 plant natural products that integrates accurate mass
27 616 – time tag to incorporate retention time relative to an internal standard in a similar fashion
28 617 as it is usually done for GC-MS based annotation, therefore, in order to use this feature it is
29 618 necessary to analyze samples with addition of the same internal standard used when
30 619 developing the database entries. It is important to highlight that this kind of approach is
31 620 much less effective for LC-MS where relative retention time is prone to larger variation. MS-
32 621 MS Fragment Viewer (<http://webs2.kazusa.or.jp/msmsfragmentviewer/>) is a very small and
33 622 not very frequently updated database containing FT-MS, IT- and FT-MS/MS spectral data on
34 623 116 flavonoids. ReSpect [172] is a collection of MSⁿ spectra data from 9,017 phytochemicals
35 624 from literature and standards with searching functionalities very similar to MassBank, and
36 625 WEIZMASS [173], a metabolite spectral library of high-resolution MS data from 3,540 plant
37 626 metabolites that uses a probabilistic approach to match library and experimental data with
38 627 the MatchWeiz software. WEIZMASS is available for implementation in R as a pipeline for
39 628 metabolite identification which can be easily integrated with data processing. While this is a
40 629 much less accessible tool for general use compared with other web based databases the
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

630 results obtained are far more considerable and the effort required in its use is, therefore,
631 more than compensated by the gains which it affords.

632 A very common issue encountered in data from mass spectrometry is the presence of a
633 variety of contaminants from sample preparation and analysis that can be challenging for
634 data interpretation. MaConDa [174] provides a very useful database of common
635 contaminants and adducts in mass spectrometry, containing over 200 contaminant records
636 with origin of the contaminant, its mass and the adducts formed. MaConDa can be
637 downloaded in different formats or accessed via the web browser.

638 Compound spectra databases are essential for identification of metabolites by mass
639 spectrometry, but a significant effort has also been directed towards the development of
640 repositories of experimental data on specific samples to facilitate dereplication studies and
641 data analysis. These databases are often restricted to specific species, as it is the case for
642 AtMetExpress [175], a LC-MS database of Arabidopsis with data on 20 different ecotypes
643 and 36 developmental stages which allows users to download raw and processed data as
644 well as query using mass chromatogram features in the web platform and visualize
645 annotation and distribution of selected features. MeKO [176], is a GC-MS database of 50
646 Arabidopsis KO mutants. All raw data can be downloaded as netCDF files and results from
647 data analysis can be visualized in a very informative summary in the web browser that
648 shows plant phenotypes, differentially accumulated metabolites indicated in a pathway map
649 and log fold changes for most significantly changed metabolites. MoTo DB [177] is a LC-MS
650 database of *Solanum lycopersicum* with information of annotated metabolites where the
651 user can search for specific masses or a range of masses. The database is based on accurate
652 mass and the user therefore does not have access to raw data and chromatograms. NaDH
653 [178], a platform for integration and visualization of different omics datasets of *Nicotiana
654 attenuata* including LC-MS data on 14 different tissues, allows search for spectra based on
655 name and m/z and provides some interesting tools for data interpretation easily accessible
656 directly from the metabolite entry including metabolite-metabolite and metabolite-gene
657 coexpression analysis and visualization of metabolite expression across different tissues in a
658 bar chart or eFP browser interface. The Optimas-DW software [179], is a data collection for
659 maize data of 15 different experiments, the interface for metabolites allows easy browsing
660 through all the metabolites and visualization of values for individual experiments in a table
661 format but no access to raw data, and the SoyMetDB [180], a metabolomics database for
662 soybean, with GC-MS and LC-MS data of four different tissues under two different
663 conditions, which has a simple interface that provide search by metabolite name or
664 browsing through the whole dataset, metabolite entries provide m/z, retention time as well
665 as an apparent defunct link to a pathway viewer. Similar databases with relative broader
666 spectra include the plant specific KOMIC Maket [181] currently warehousing LC-MS data on
667 74 samples from 17 species, in which the user can search for peaks and browse through
668 samples and the interface shows retention times, m/z and annotation details classifying the
669 annotation based on a grading system. MS2T [182] is an MSMS library created using a

670 function for automatic Tandem MS acquisition from over 150 samples from 10 different
1 671 plant species, the web platforms allows search by retention time, m/z and spectra similarity.
2
3 672 PMR [183], is a database for plants and eukaryotic microorganisms which includes the
4
5 673 earlier database of medicinal plants MPMR [184] and currently comprises of GC-MS and LC-
6
7 674 MS data on 24 species from different sources and experiments including different tissues
8
9 675 and developmental stages. It has an easy and clear interface with summary of all the
10
11 676 experiments once an individual species is selected including metadata and annotated
12
13 677 metabolites. It additionally allows the download of all the results in csv format in the form
14
15 678 of peak tables and it has some basic tool for comparative analysis where volcano plots can
16
17 679 be generated comparing different experiments. By contrast, the more general databases
18
19 680 Bio-MassBank (<http://bio.massbank.jp/>), a repository of LC-MS and GC-MS data from
20
21 681 biological samples, in contrast with the original MassBank in this database most of the data
22
23 682 is tagged as “Unknowns” or are just putative metabolites, searching functions are similar to
24
25 683 the original database but it includes a samples section where it is possible to access all the
26
27 684 experiments available. MassBase (<http://webs2.kazusa.or.jp/massbase/>) is a large
28
29 685 repository providing raw and processed mass chromatograms on 46,398 samples of over 40
30
31 686 species, including several plants, analyzed by LC-MS, GC-MS and CE-MS. Metabolomics
32
33 687 Workbench [185] is a repository of a variety of metabolomics experiments containing over
34
35 688 60,000 entries, including raw and processed MS data, a section with detailed protocols for
36
37 689 the experiments, and web tools for analysis and interpretation that can be used with any
38
39 690 uploaded data. Similarly, Metabolights [186], is a cross species repository containing data
40
41 691 from 190 mass spectrometry based metabolomics studies that is currently recommended as
42
43 692 repository of experimental data by many journals, all experimental data can be downloaded
44
45 693 from an ftp server and data submission is powered by the use of ISA software that assists in
46
47 694 the reporting and management of metadata. MetabolomeXchange [187], is a data
48
49 695 aggregation system that allows users to efficiently explore experimental metabolomics data
50
51 696 from different databases including MetaboLights and Metabolomics Workbench providing
52
53 697 an RSS feeding service to allow users to get updates over the datasets available. Similarly,
54
55 698 GNPS [188], a plant natural product knowledge base for community-wide organization and
56
57 699 sharing of raw, processed or identified tandem mass spectrometry data currently
58
59 700 comprising of 221,083 MS/MS spectra from 18,163 unique compounds. The platform allows
60
61 701 users to upload data and provides a series of tools for analysis and interpretation based on
62
63 702 the data from the database.

50 703 As previously mentioned, many resources that are particularly useful for data interpretation
51 704 organize the data in pathways based on literature data, and often also provide tools for data
52 705 visualization and interpretation. Many of these databases contain either generic pathways
53 706 or combine different organisms, some examples are KEGG [189], which includes 504
54 707 pathway maps with 17,891 compounds and 10,419 reactions for 4,607 different organisms,
55 708 representing data in an interactive interface that links the entries to a great amount of
56 709 external resources being one of the most popular sources of information on metabolic
60
61
62
63
64
65

710 pathways One of the greatest issues of KEGG leading many user to misinterpreting their
1 711 data is that it displays all genes in generic pathway maps of which some are characterized
2 712 only by similarity, resulting in pathways that are not present in the analysed organism being
3 713 represented. By contrast, WikiPathways [190], is a wiki-style website with 2,471 community
4 714 curated pathways of 28 different organisms. Its interactive interface is similar to KEGG
5 715 providing link with external resources for metabolites and enzymes. Similarly, kpath [191], is
6 716 a database that integrates information related to metabolic pathways with 74,180 pathways
7 717 13,153 reactions and 37,029 metabolites providing tools for pathway visualization, editing
8 718 and relationship search. BioCyc [192], is a collection of 9,387 Pathway/Genome Databases,
9 719 and MetaCyc [192] is the largest curated database of experimentally elucidated metabolic
10 720 pathways containing 2,491 pathways from 2,816 different organisms. KBase [193],
11 721 meanwhile, is a data platform with data on plants and microbes that allow users to upload
12 722 their own data and integrates data and tools for systems biology including 1,470 metabolic
13 723 pathways with 33,773 reactions and 27,838 compounds, genome data on 60 different plant
14 724 species and tools for assembly, annotation, metabolic modeling, comparative analysis,
15 725 phylogenetic analysis and expression analysis. There are also a significant amount of plant
16 726 specific data organized in databases like KaPPA-View4 [141], containing 153 pathways with
17 727 1,427 compounds and 1,434 reaction from 10 species, allowing users to upload their own
18 728 data and is able to represent gene-to-gene and metabolite-to-metabolite relationships as
19 729 curves on a metabolic pathway maps to help in data interpretation. PlantCyc
20 730 (<http://www.plantcyc.org/>) provides access to manually curated or reviewed information
21 731 about metabolic pathways in over 800 pathways of 350 plant species, usefully the platform
22 732 provides “evidence codes” to clearly indicate the type of support associated with each
23 733 database item. MetaCrop [194], is a pathway database containing information about seven
24 734 major crop plants and two model plants that allows integration of experimental data into
25 735 metabolic pathways, as well as the automatic export of information for the creation of
26 736 detailed metabolic models. Similarly, MetNetDB [195], contains integrative information on
27 737 metabolic and regulatory networks of Arabidopsis and Soybean with metabolism, signalling,
28 738 and transcriptional pathways being fully integrated into a single network and manually
29 739 curated subcellular localization is represented in the pathway maps. The network
30 740 information can be exported to other applications for network analysis such as exploRase,
31 741 and Cytoscape/FCM. Like MetNetDB, Gramene [196] is an integrated data resource for
32 742 comparative functional genomics in crops and model plants that host pathway databases for
33 743 rice, maize, Bracypodium, and sorghum as well as providing mirrors for MetaCyc and
34 744 PlantCyc data. It is worth mentioning a few resources that are focused on the reactions
35 745 within the pathways offering detailed curated metabolic reactions, namely BioMeta [197],
36 746 whose contents are based on the KEGG Ligand database with a large number of chemical
37 747 structures corrected with respect to constitution and reactions’ stereochemistry being
38 748 correctly balanced. BKM-react [198] is a non-redundant biochemical reaction database
39 749 containing 18,172 unique biochemical reactions retrieved from BRENDA, KEGG, and
40 750 MetaCyc databases that were matched and integrated by aligning substrates and products.

751 Similar to this MetRxn [199], also integrates information from BRENDA, KEGG and MetaCyc,
752 combining also Reactome.org and 44 metabolic models in a standardized description of
753 metabolites and reactions where all metabolites have matched synonyms, resolved
754 protonation states, and are linked to unique structures, and all reactions are balanced.

755 Together with the development of many prediction tools previously mentioned we watched
756 in the last years the development of some interesting *In Silico* databases that are extremely
757 useful for *de novo* metabolite identification such as MINE [200], a database developed by
758 the integration of an algorithm called Biochemical Network Integrated Computational
759 Explorer (BNICE) and expert-curated reaction rules to predict chemical structures product of
760 enzyme promiscuity, MetCCS [201] a database and algorithm for prediction of Collision
761 Cross-Section values for metabolites in ion mobility mass spectrometry, a technique
762 increasingly used to assist metabolite elucidation based on the drift speed of the ion that is
763 proportional to its cross section, and the plant specific ISDB [202] an *in silico* database of
764 natural products generated using CFM-ID [126] with input from the commercial Dictionary
765 of Natural Products.

766 **Other programs of interest**

767 The complexity of metabolomics data experiments, particularly in terms of sample number
768 and metadata pushed the development of many tools for experiment and metadata
769 management, and while many of these functions are integrated in some of the databases
770 previously discussed there are a few specialized tools such as QTREDS [203] and MASTR-MS
771 [204], that are LIMS based software for assisting in organizing experimental design,
772 metadata management and sample data acquisition , MetaDB [205] a web application for
773 Metabolomics metadata management with interface to MetaMS data processing tool, and
774 Metabolonote [206], a metadata database/management system.

775 The enormous amount of data available for metabolomics raises many questions regarding
776 how to easily access and unify all this data, taking into account the vast chemical space
777 explored in these experiments. Many tools have been developed with the purpose of
778 facilitating access to chemical data spread in the literature, from the development of
779 identifiers to reduce duplication of information such as the SPLASH [207] hash designed for
780 the MoNA database, to tools like Metmask [208], for managing different identifiers, CTS
781 [209], for translation of chemical identifiers, PhenoMeter [210] for querying databases
782 based on metabolic phenotype and Metab2MeSH [211] for a more efficient literature
783 search that automatically annotate compounds with the concepts defined in MeSH
784 providing a fast link between compound and the literature.

785 Different vendors usually export their data in proprietary formats which complicates data
786 transfer across different platforms. Most proprietary software are able to convert files to
787 .cdf format, but some tools from which the most popular is msConverter from Proteowizard
788 (<http://proteowizard.sourceforge.net/>) can handle conversion from/to different formats

789 including mzXML. mzTab is another format proposed by the Proteomics Standards Initiative
790 targeting researchers outside of proteomics, it is supposed to contain the minimal
791 information required to evaluate the results of a proteomics experiment making it more
792 accessible to non-experts, jmzTab [212] is a java application that provides reading and
793 writing capabilities and conversion of files to mzTab. The PeakML [213] file format is an
794 initiative developed by the creators of mzMatch to enable the exchange of data between
795 analysis software by representing peak and meta-information from each step in an analysis
796 pipeline, as a proof of concept the R-package 'mzmatch.R' was developed to extend XCMS
797 functionalities for storing and reading data in PeakML format.

798 All equipment for mass spectrometry comes with their own software for data visualization
799 and some basic analysis but those are usually not designed to deal with the complexities of
800 metabolomics datasets. There are some interesting open source alternatives such as
801 BatMass [214] and Mass++ [215] for data visualization, and for generating images from raw
802 data like SpeckTackle [216] that provides several pre-defined chart types easy to integrate
803 into web-facing resources and RMassBank [217] capable of automatically generating
804 MassBank records from raw MS and MS/MS data.

805 Mass spectrometry imaging is a relative young technique that has been growing fast in
806 importance providing high resolution spatial distribution of small molecules in molecular
807 histology [218]. Few tools have been developed so far, namely EXIMS [219] for data
808 processing and analysis, and OpenMSI [220], a web-based visualization, analysis and
809 management tool.

810 Lipidomics data requires a very specialized pipeline and therefore many tools were
811 developed exclusively for this kind of analysis however we will only briefly summarize these
812 here. ALEX [221], MRM-DIFF [222], LICRE [223], LipidXplorer [224], LIMSA [225], VaLID
813 [226], LOBSTAHS [227], Lipid-Pro [228], LDA [229] and LipidQA [230] are all tools for
814 processing, annotating and analyzing lipidomics data. Lipids databases include LIPID MAPS
815 [231], LIPIDBANK [232], LipidBlast [233], and in silico generated lipids database LipidHome
816 [234], SwissLipids [235] and ARALIP
817 (<http://aralip.plantbiology.msu.edu/pathways/pathways>).

818 **Future perspectives**

819 Many of the resources presented here were fruit of the efforts in setting the theoretical
820 background for each step in the data processing and analysis workflow. However, more
821 recent efforts are moving towards the development of integrated tools, which are often
822 developed by the integration of already well established tools into a single pipeline in an
823 attempt to accelerate the process and in a few cases providing an easier interface. XCMS
824 online, for example, is a web platform providing most of the function from XCMS with
825 additional capabilities for interactive exploratory data visualization and analysis in a much
826 easier interface than the original software [236], HayStack [237], is a web platform that uses

827 XCMS to process data and automatically generates total ion chromatograms (TIC) and base
1 828 peak chromatograms as well as offering an easy way of plotting extracted ion
2 chromatograms (EIC) and some basic statistical tools such as PCA scores plot, volcano plots,
3 829 and dendrograms for group comparisons, SMART [238] is an R package that combines
4 830 different tools such as XCMS and CAMERA with a series of common statistical approaches to
5 831 provide an integrated pipeline for data processing, visualization, and analysis. MZmine 2
6 832 [239] is another very popular tool with over 1000 citations, it was originally developed for
7 833 LC-MS data processing but it became one of the most popular platforms for development of
8 834 integrated tools in Java providing a user-friendly, flexible and extendable software
9 835 constantly updated and with a set of modules covering most steps of LC-MS processing and
10 836 data analysis workflow including several option of visualization tools. MetSign [240] is a
11 837 MATLAB package providing tools for spectra deconvolution, metabolite putative assignment
12 838 by matching m/z and peak isotopic distribution against its own database, peak list
13 839 alignment, a series of normalization algorithms, statistical significance tests, unsupervised
14 840 clustering, and time course analysis, all in a modular and interactive design presented with a
15 841 wizard to facilitate the analysis workflow. MultiAlign [241] is a software developed in the
16 842 .NET platform using C++ and C# originally for proteomics but that can also be used for
17 843 metabolomics comparative analysis, its functionalities include feature detection, alignment,
18 844 several plotting options, normalization, and basic statistical comparisons, Metabolome
19 845 Express [242] works as a web server to process, interpret and share GC/MS metabolomics
20 846 datasets, whilst MAIT [243] is an R package aiming at providing an end-to-end
21 847 programmable metabolomics pipeline with emphasis in metabolite annotation and
22 848 statistics, it uses XCMS for peak detection, an approach based on CAMERA combined with
23 849 an user defined table of biotransformations followed by database search for metabolite
24 850 annotation and a series of statistical tests to identify statistically significant features
25 851 containing the highest amount of class-related information. By contrast, MAVEN [244] is a
26 852 software for data processing, analysis and visualization with some interesting features for
27 853 pathway-based visualization of isotope-labeling data that can be helpful for the
28 854 interpretation of this kind of experiment. MeltDB [245] is a java web based platform that
29 855 integrates different algorithms for data processing, compound identification by spectra
30 856 matching statistical analysis, data visualization and integration with transcriptomics and
31 857 proteomics datasets via the ProMeTra software. It provides a tool for saving peaks of
32 858 reference compounds directly in the MeltDB database, and allows storage and sharing of
33 859 projects within the web server. MetaboAnalyst [246] is another java web platform with data
34 860 processing and a comprehensive set of data analysis tools, it includes most common
35 861 approaches for statistical analysis as well as modules for functional enrichment analysis,
36 862 metabolic pathway analysis, time series and two-factor data analysis, biomarker analysis,
37 863 sample size and power analysis, integrated pathway analysis, and image and report
38 864 generation. The program mzMatch [213] is a popular Java toolkit for processing, filtering,
39 865 and annotation, with particular focus on integration of processed data across different
40 866 platforms and providing a customizable modular pipeline to facilitate the development and
41 867

868 integration of different tools. It includes many other tools previously described here like
1 869 mzmatchISO and metAssign and it is based entirely in the PeakML file format. The MarVis-
2 870 Suite [247] is a software for the interactive ranking, filtering, combination, clustering,
3 871 visualization, and functional analysis of transcriptomics and metabolomics data sets, the
4 872 clustering algorithm is based on one-dimensional self-organizing maps (1D-SOMs), and the
5 873 software additionally provides functions for metabolite annotation and pathway
6 874 reconstruction. MetMSLine [248] is an R package that works with processed data providing
7 875 a series of statistical analysis steps focusing on biomarker discovery combined with
8 876 metabolite annotation based on exact mass matching against a target list of metabolites
9 877 and MassCascade [249] is a Java library that takes advantage of the KINIME workflow
10 878 environment facilitating integration with other tools and making the tool user friendly, the
11 879 core library contains a collection of data processing algorithms, a visualization framework
12 880 and metabolite annotation functions, while the plug-in for KNIME allows easy integration
13 881 with other statistical workflows. MASSyPup [250] does not actually integrate different
14 882 procedures but it does provides an easy platform for accessing many different tools in the
15 883 form of a Linux distribution that can be run directly from different media without
16 884 installation.

885 It is clear from this review the infinity of choices for performing a variety of functions and
886 the fast pace by which they change and get outdated; hence it is an arduous task to keep
887 updated of all of them. Some research groups, engaged in the development of
888 metabolomics tools, have their own repositories like KOMICS [251], MetaOpen
889 (<http://metaopen.sourceforge.net/>) and PRIME [252], while OMICtools [253], NAR online
890 Molecular Biology Database Collection and the Bioinformatics Links Directory provide
891 unified repositories but still covering only a small portion of all the resources available. With
892 the rapid development of new tools it is of great interest for the metabolomics community
893 to develop classification systems and repositories to catalog and provide a platform for
894 submission, curation and feedback facilitating users' access to the most appropriate
895 resources for each aim. Another clear observation that can be made from the proceeding
896 sections is that the number of tools for analysis by far exceeds that of the number of data
897 repositories whilst metabolomics is clearly difficult to fully standardize this is still a great
898 shame. There are a number of clear reporting standards that should aid in this respect [254],
899 furthermore, both the existing databases and carefully compared meta-analysis [22, 255],
900 demonstrate that such approaches are indeed highly powerful in the enhancement of
901 biological understanding. As such we feel that it is an urgent priority to focus efforts on the
902 improvement of this feature of computational metabolomics since it will aid not only in the
903 expansion of our coverage of the metabolite complement of the plant cell but also in the
904 equally important task of interpreting the biological function of the individual metabolites
905 themselves.

906 **Competing interests**

907 The authors declare that they have no competing interests.

908

909 **Acknowledgement**

910 We thank the Max Planck Society, the National Council for Scientific and Technological
911 Development CNPq-Brazil (LPS) and the IMPRS-PMPG program (TN) for the financial
912 support.

913 **References**

- 914 1. Oliver SG, Winson MK, Kell DB and Baganz F. Systematic functional analysis of the yeast
915 genome. *Trends Biotechnol.* 1998;16 9:373-8. doi:10.1016/s0167-7799(98)01214-1.
- 916 2. Fiehn O, Kopka J, Dormann P, Altmann T, Trethewey RN and Willmitzer L. Metabolite
917 profiling for plant functional genomics. *Nat Biotechnol.* 2000;18 11:1157-61.
918 doi:10.1038/81137.
- 919 3. Sauter H, Lauer M and Fritsch H. METABOLIC PROFILING OF PLANTS - A NEW DIAGNOSTIC-
920 TECHNIQUE. *Abstr Pap Am Chem Soc.* 1988;195:129-AGRO.
- 921 4. Dorr JR, Yu Y, Milanovic M, Beuster G, Zasada C, Dabritz JHM, et al. Synthetic lethal
922 metabolic targeting of cellular senescence in cancer therapy. *Nature.* 2013;501 7467:421-+.
923 doi:10.1038/nature12437.
- 924 5. Kell DB. Metabolomics and systems biology: making sense of the soup. *Current Opinion in*
925 *Microbiology.* 2004;7 3:296-307. doi:10.1016/j.mib.2004.04.012.
- 926 6. Nicholson JK and Wilson ID. Understanding 'global' systems biology: Metabonomics and the
927 continuum of metabolism. *Nature Reviews Drug Discovery.* 2003;2 8:668-76.
928 doi:10.1038/nrd1157.
- 929 7. Fernie AR and Schauer N. Metabolomics-assisted breeding: a viable option for crop
930 improvement? *Trends in Genetics.* 2009;25 1:39-48. doi:10.1016/j.tig.2008.10.010.
- 931 8. Meyer RC, Steinfath M, Lisek J, Becher M, Witucka-Wall H, Torjek O, et al. The metabolic
932 signature related to high plant growth rate in *Arabidopsis thaliana*. *Proceedings of the*
933 *National Academy of Sciences of the United States of America.* 2007;104 11:4759-64.
934 doi:10.1073/pnas.0609709104.
- 935 9. Roessner U, Willmitzer L and Fernie AR. Metabolic profiling and biochemical phenotyping of
936 plant systems. *Plant Cell Reports.* 2002;21 3:189-96. doi:10.1007/s00299-002-0510-8.
- 937 10. Schauer N and Fernie AR. Plant metabolomics: towards biological function and mechanism.
938 *Trends in Plant Science.* 2006;11 10:508-16. doi:10.1016/j.tplants.2006.08.007.
- 939 11. Weckwerth W. Metabolomics in systems biology. *Annu Rev Plant Biol.* 2003;54:669-89.
940 doi:10.1146/annurev.arplant.54.031902.135014.
- 941 12. Fernie AR and Stitt M. On the Discordance of Metabolomics with Proteomics and
942 Transcriptomics: Coping with Increasing Complexity in Logic, Chemistry, and Network
943 Interactions. *Plant Physiology.* 2012;158 3:1139-45. doi:10.1104/pp.112.193235.
- 944 13. Nobeli I, Ponstingl H, Krissinel EB and Thornton JM. A structure-based anatomy of the E-coli
945 metabolome. *Journal of Molecular Biology.* 2003;334 4:697-719.
946 doi:10.1016/j.jmb.2003.10.008.
- 947 14. van der Werf MJ, Overkamp KM, Mulwijk B, Coulier L and Hankemeier T. Microbial
948 metabolomics: Toward a platform with full metabolome coverage. *Analytical Biochemistry.*
949 2007;370 1:17-25. doi:10.1016/j.ab.2007.07.022.
- 950 15. Tohge T, Scossa F and Fernie AR. Integrative Approaches to Enhance Understanding of Plant
951 Metabolic Pathway Structure and Regulation. *Plant Physiology.* 2015;169 3:1499-511.
952 doi:10.1104/pp.15.01006.

953 16. Sulpice R, Pyl E-T, Ishihara H, Trenkamp S, Steinfath M, Witucka-Wall H, et al. Starch as a
1 954 major integrator in the regulation of plant growth. *Proceedings of the National Academy of*
2 955 *Sciences*. 2009;106 25:10348-53. doi:10.1073/pnas.0903478106.

3 956 17. Davey MP, Burrell MM, Woodward FI and Quick WP. Population-specific metabolic
4 957 phenotypes of *Arabidopsis lyrata* ssp. *petraea*. *New Phytologist*. 2008;177 2:380-8.
5 958 doi:10.1111/j.1469-8137.2007.02282.x.

6 959 18. Beleggia R, Rau D, Laidò G, Platani C, Nigro F, Fragasso M, et al. Evolutionary Metabolomics
7 960 Reveals Domestication-Associated Changes in Tetraploid Wheat Kernels. *Molecular Biology*
8 961 *and Evolution*. 2016;33 7:1740-53. doi:10.1093/molbev/msw050.

9 962 19. Kliebenstein D. Advancing Genetic Theory and Application by Metabolic Quantitative Trait
10 963 Loci Analysis. *The Plant Cell*. 2009;21 6:1637-46. doi:10.1105/tpc.109.067611.

11 964 20. Luo J. Metabolite-based genome-wide association studies in plants. *Current Opinion in Plant*
12 965 *Biology*. 2015;24:31-8. doi:<http://dx.doi.org/10.1016/j.pbi.2015.01.006>.

13 966 21. Brotman Y, Landau U, Pnini S, Lisek J, Balazadeh S, Mueller-Roeber B, et al. The LysM
14 967 receptor-like kinase LysM RLK1 is required to activate defense and abiotic-stress responses
15 968 induced by overexpression of fungal chitinases in *Arabidopsis* plants. *Molecular plant*.
16 969 2012;5 5:1113-24.

17 970 22. Obata T and Fernie AR. The use of metabolomics to dissect plant responses to abiotic
18 971 stresses. *Cellular and Molecular Life Sciences*. 2012;69 19:3225-43. doi:10.1007/s00018-012-
19 972 1091-5.

20 973 23. Tohge T and Fernie AR. Web-based resources for mass-spectrometry-based metabolomics: A
21 974 user's guide. *Phytochemistry*. 2009;70 4:450-6.
22 975 doi:<http://dx.doi.org/10.1016/j.phytochem.2009.02.004>.

23 976 24. Hibbert DB. Experimental design in chromatography: A tutorial review. *Journal of*
24 977 *Chromatography B*. 2012;910:2-13. doi:<http://dx.doi.org/10.1016/j.jchromb.2012.01.020>.

25 978 25. Gullberg J, Jonsson P, Nordström A, Sjöström M and Moritz T. Design of experiments: an
26 979 efficient strategy to identify factors influencing extraction and derivatization of *Arabidopsis*
27 980 *thaliana* samples in metabolomic studies with gas chromatography/mass spectrometry.
28 981 *Analytical Biochemistry*. 2004;331 2:283-95. doi:<http://dx.doi.org/10.1016/j.ab.2004.04.037>.

29 982 26. Nistor I, Cao M, Debrus B, Lebrun P, Lecomte F, Rozet E, et al. Application of a new
30 983 optimization strategy for the separation of tertiary alkaloids extracted from *Strychnos*
31 984 *usambarensis* leaves. *Journal of Pharmaceutical and Biomedical Analysis*. 2011;56 1:30-7.
32 985 doi:<http://dx.doi.org/10.1016/j.jpba.2011.04.027>.

33 986 27. Bradbury J, Genta-Jouve G, Allwood JW, Dunn WB, Goodacre R, Knowles JD, et al. MUSCLE:
34 987 automated multi-objective evolutionary optimization of targeted LC-MS/MS analysis.
35 988 *Bioinformatics*. 2015;31 6:975-7. doi:10.1093/bioinformatics/btu740.

36 989 28. Nikolskiy I, Siuzdak G and Patti GJ. Discriminating precursors of common fragments for large-
37 990 scale metabolite profiling by triple quadrupole mass spectrometry. *Bioinformatics*. 2015;31
38 991 12:2017-23.

39 992 29. Katajamaa M and Orešič M. Data processing for mass spectrometry-based metabolomics.
40 993 *Journal of Chromatography A*. 2007;1158 1-2:318-28.
41 994 doi:<http://dx.doi.org/10.1016/j.chroma.2007.04.021>.

42 995 30. Sugimoto M, Kawakami M, Robert M, Soga T and Tomita M. Bioinformatics tools for mass
43 996 spectroscopy-based metabolomic data processing and analysis. *Current bioinformatics*.
44 997 2012;7 1:96-108.

45 998 31. Lange E, Tautenhahn R, Neumann S and Gröpl C. Critical assessment of alignment
46 999 procedures for LC-MS proteomics and metabolomics measurements. *BMC Bioinformatics*.
50 1000 2008;9:375-. doi:10.1186/1471-2105-9-375.

51 1001 32. Tautenhahn R, Böttcher C and Neumann S. Highly sensitive feature detection for high
52 1002 resolution LC/MS. *BMC Bioinformatics*. 2008;9 1:504. doi:10.1186/1471-2105-9-504.

- 1003 33. Lommen A. MetAlign: interface-driven, versatile metabolomics tool for hyphenated full-scan
1 1004 mass spectrometry data preprocessing. *Analytical Chemistry*. 2009;81 8:3079-86.
2 1005 34. Smith CA, Want EJ, O'Maille G, Abagyan R and Siuzdak G. XCMS: processing mass
3 1006 spectrometry data for metabolite profiling using nonlinear peak alignment, matching, and
4 1007 identification. *Analytical Chemistry*. 2006;78 doi:10.1021/ac051437y.
5 1008 35. Tengstrand E, Lindberg J and Åberg KM. TracMass 2 A Modular Suite of Tools for Processing
6 1009 Chromatography-Full Scan Mass Spectrometry Data. *Analytical Chemistry*. 2014;86 7:3435-
7 1010 42.
8 1011 36. Chang H-Y, Chen C-T, Lih TM, Lynn K-S, Juo C-G, Hsu W-L, et al. iMet-Q: A User-Friendly Tool
9 1012 for Label-Free Metabolomics Quantitation Using Dynamic Peak-Width Determination. *PLoS*
10 1013 *One*. 2016;11 1:e0146112. doi:10.1371/journal.pone.0146112.
11 1014 37. Treviño V, Yañez-Garza IL, Rodríguez-López CE, Urrea-López R, Garza-Rodríguez ML, Barrera-
12 1015 Saldaña HA, et al. GridMass: a fast two-dimensional feature detection method for LC/MS.
13 1016 *Journal of Mass Spectrometry*. 2015;50 1:165-74.
14 1017 38. Duran AL, Yang J, Wang L and Sumner LW. Metabolomics spectral formatting, alignment and
15 1018 conversion tools (MSFACTS). *Bioinformatics*. 2003;19 17:2283-93.
16 1019 39. Broeckling CD, Reddy IR, Duran AL, Zhao X and Sumner LW. MET-IDEA: data extraction tool
17 1020 for mass spectrometry-based metabolomics. *Analytical Chemistry*. 2006;78 13:4334-41.
18 1021 40. Fructuoso S, Sevilla Á, Bernal C, Lozano AB, Iborra JL and Cánovas M. EasyLCMS: an
19 1022 asynchronous web application for the automated quantification of LC-MS data. *BMC*
20 1023 *research notes*. 2012;5 1:428.
21 1024 41. Creek DJ, Jankevics A, Burgess KE, Breitling R and Barrett MP. IDEOM: an Excel interface for
22 1025 analysis of LC-MS-based metabolomics data. *Bioinformatics*. 2012;28 7:1048-9.
23 1026 42. Conley CJ, Smith R, Torgrip RJ, Taylor RM, Tautenhahn R and Prince JT. Massifquant: open-
24 1027 source Kalman filter-based XC-MS isotope trace feature detection. *Bioinformatics*. 2014;30
25 1028 18:2636-43.
26 1029 43. Zhang W, Chang J, Lei Z, Huhman D, Sumner LW and Zhao PX. MET-COFEA: a liquid
27 1030 chromatography/mass spectrometry data processing platform for metabolite compound
28 1031 feature extraction and annotation. *Analytical Chemistry*. 2014;86 13:6245-53.
29 1032 44. Zhang W, Lei Z, Huhman D, Sumner LW and Zhao PX. MET-XAlign: A metabolite cross-
30 1033 alignment tool for LC/MS-based comparative metabolomics. *Analytical Chemistry*. 2015;87
31 1034 18:9114-9.
32 1035 45. Yu T, Park Y, Johnson JM and Jones DP. apLCMS—adaptive processing of high-resolution
33 1036 LC/MS data. *Bioinformatics*. 2009;25 15:1930-6.
34 1037 46. Uppal K, Soltow QA, Strobel FH, Pittard WS, Gernert KM, Yu T, et al. xMSanalyzer: automated
35 1038 pipeline for improved feature detection and downstream analysis of large-scale, non-
36 1039 targeted metabolomics data. *BMC Bioinformatics*. 2013;14 1:15.
37 1040 47. Myint L, Kleensang A, Zhao L, Hartung T and Hansen KD. Joint bounding of peaks across
38 1041 samples improves differential analysis in mass spectrometry-based metabolomics. *Analytical*
39 1042 *Chemistry*. 2017; doi:10.1021/acs.analchem.6b04719.
40 1043 48. Wandy J, Daly R, Breitling R and Rogers S. Incorporating peak grouping information for
41 1044 alignment of multiple liquid chromatography-mass spectrometry datasets. *Bioinformatics*.
42 1045 2015;31 12:1999-2006.
43 1046 49. Wehrens R, Bloemberg TG and Eilers PH. Fast parametric time warping of peak lists.
44 1047 *Bioinformatics*. 2015;31 18:3063-5.
45 1048 50. Stein SE. An integrated method for spectrum extraction and compound identification from
46 1049 gas chromatography/mass spectrometry data. *Journal of the American Society for Mass*
47 1050 *Spectrometry*. 1999;10 8:770-81. doi:[http://dx.doi.org/10.1016/S1044-0305\(99\)00047-1](http://dx.doi.org/10.1016/S1044-0305(99)00047-1).
48 1051 51. Aggio R, Villas SG and Ruggiero K. Metab: an R package for high-throughput analysis of
49 1052 metabolomics data generated by GC-MS. *Bioinformatics*. 2011;27 16:2316-8.
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1053 52. Bunk B, Kucklick M, Jonas R, Münch R, Schobert M, Jahn D, et al. MetaQuant: a tool for the
1 1054 automatic quantification of GC/MS-based metabolome data. *Bioinformatics*. 2006;22
2 1055 23:2962-5.
- 3 1056 53. Hiller K, Hangebrauk J, Jäger C, Spura J, Schreiber K and Schomburg D. MetaboliteDetector:
4 1057 comprehensive analysis tool for targeted and nontargeted GC/MS based metabolome
5 1058 analysis. *Analytical Chemistry*. 2009;81 9:3429-39.
- 7 1059 54. Luedemann A, Strassburg K, Erban A and Kopka J. TagFinder for the quantitative analysis of
8 1060 gas chromatography—mass spectrometry (GC-MS)-based metabolite profiling experiments.
9 1061 *Bioinformatics*. 2008;24 5:732-7.
- 10 1062 55. Cuadros-Inostroza Á, Caldana C, Redestig H, Kusano M, Lisec J, Peña-Cortés H, et al.
11 1063 TargetSearch—a Bioconductor package for the efficient preprocessing of GC-MS metabolite
12 1064 profiling data. *BMC Bioinformatics*. 2009;10 1:428.
- 14 1065 56. O'Callaghan S, De Souza DP, Isaac A, Wang Q, Hodkinson L, Olshansky M, et al. PyMS: a
15 1066 Python toolkit for processing of gas chromatography-mass spectrometry (GC-MS) data.
16 1067 Application and comparative study of selected tools. *BMC Bioinformatics*. 2012;13 1:115.
- 17 1068 57. Jellema RH, Krishnan S, Hendriks MM, Muilwijk B and Vogels JT. Deconvolution using signal
18 1069 segmentation. *Chemometrics and Intelligent Laboratory Systems*. 2010;104 1:132-9.
- 20 1070 58. Wehrens R, Weingart G and Mattivi F. metaMS: An open-source pipeline for GC–MS-based
21 1071 untargeted metabolomics. *Journal of Chromatography B*. 2014;966:109-16.
- 22 1072 59. Kuich PHJ, Hoffmann N and Kempa S. Maui-VIA: a user-friendly software for visual
23 1073 identification, alignment, correction, and quantification of gas chromatography—mass
24 1074 spectrometry data. *Frontiers in bioengineering and biotechnology*. 2014;2.
- 26 1075 60. Domingo-Almenara X, Brezmes J, Vinaixa M, Samino S, Ramirez N, Ramon-Krauel M, et al.
27 1076 eRah: A Computational Tool Integrating Spectral Deconvolution and Alignment with
28 1077 Quantification and Identification of Metabolites in GC/MS-Based Metabolomics. *Analytical
29 1078 Chemistry*. 2016;88 19:9821-9.
- 31 1079 61. Ni Y, Su M, Qiu Y, Jia W and Du X. ADAP-GC 3.0: Improved Peak Detection and Deconvolution
32 1080 of Co-eluting Metabolites from GC/TOF-MS Data for Metabolomics Studies. *Analytical
33 1081 Chemistry*. 2016;88 17:8802-11.
- 34 1082 62. Wei X, Shi X, Koo I, Kim S, Schmidt RH, Arteel GE, et al. MetPP: a computational platform for
35 1083 comprehensive two-dimensional gas chromatography time-of-flight mass spectrometry-
36 1084 based metabolomics. *Bioinformatics*. 2013;29 14:1786-92.
37 1085 doi:10.1093/bioinformatics/btt275.
- 39 1086 63. Kuhl C, Tautenhahn R, Böttcher C, Larson TR and Neumann S. CAMERA: An Integrated
40 1087 Strategy for Compound Spectra Extraction and Annotation of Liquid Chromatography/Mass
41 1088 Spectrometry Data Sets. *Analytical Chemistry*. 2012;84 1:283-9. doi:10.1021/ac202450g.
- 42 1089 64. Alonso A, Julià A, Beltran A, Vinaixa M, Díaz M, Ibañez L, et al. AStream: an R package for
43 1090 annotating LC/MS metabolomic data. *Bioinformatics*. 2011;27 9:1339-40.
- 45 1091 65. Kessler N, Walter F, Persicke M, Albaum SP, Kalinowski J, Goesmann A, et al. Allocator: An
46 1092 interactive web platform for the analysis of metabolomic LC-ESI-MS datasets, enabling semi-
47 1093 automated, user-revised compound annotation and mass isotopomer ratio analysis. *PLoS
48 1094 One*. 2014;9 11:e113909.
- 49 1095 66. Tikunov Y, Laptinok S, Hall R, Bovy A and De Vos R. MSClust: a tool for unsupervised mass
50 1096 spectra extraction of chromatography-mass spectrometry ion-wise aligned data.
51 1097 *Metabolomics*. 2012;8 4:714-8.
- 53 1098 67. Broeckling CD, Afsar F, Neumann S, Ben-Hur A and Prenni J. RAMClust: a novel feature
54 1099 clustering method enables spectral-matching-based annotation for metabolomics data.
55 1100 *Analytical Chemistry*. 2014;86 14:6812-7.
- 56 1101 68. Gu H, Gowda GN, Neto FC, Opp MR and Raftery D. RAMSY: ratio analysis of mass
57 1102 spectrometry to improve compound identification. *Analytical Chemistry*. 2013;85 22:10771-
58 1103 9.
- 59 1103
60
61
62
63
64
65

- 1104 69. Chen G, Cui L, Teo GS, Ong CN, Tan CS and Choi H. MetTailor: dynamic block summary and
1 1105 intensity normalization for robust analysis of mass spectrometry data in metabolomics.
2 1106 *Bioinformatics*. 2015:btv434.
- 3 1107 70. Chawade A, Alexandersson E and Levander F. Normalyzer: a tool for rapid evaluation of
4 1108 normalization methods for omics data sets. *Journal of Proteome Research*. 2014;13 6:3114-
5 1109 20.
- 6 1110 71. Fernández-Albert F, Llorach R, Garcia-Aloy M, Ziyatdinov A, Andres-Lacueva C and Perera A.
7 1111 Intensity drift removal in LC/MS metabolomics by common variance compensation.
8 1112 *Bioinformatics*. 2014;30 20:2899-905.
- 9 1113 72. Shen X, Gong X, Cai Y, Guo Y, Tu J, Li H, et al. Normalization and integration of large-scale
10 1114 metabolomics data using support vector regression. *Metabolomics*. 2016;12 5:1-12.
- 11 1115 73. Karpievitch YV, Nikolic SB, Wilson R, Sharman JE and Edwards LM. Metabolomics Data
12 1116 Normalization with EigenMS. *PLoS One*. 2015;9 12:e116221.
13 1117 doi:10.1371/journal.pone.0116221.
- 14 1118 74. Huege J, Goetze J, Dethloff F, Junker B and Kopka J. Quantification of stable isotope label in
15 1119 metabolites via mass spectrometry. *Plant Chemical Genomics: Methods and Protocols*.
16 1120 2014:213-23.
- 17 1121 75. Millard P, Letisse F, Sokol S and Portais J-C. IsoCor: correcting MS data in isotope labeling
18 1122 experiments. *Bioinformatics*. 2012;28 9:1294-6.
- 19 1123 76. Jungreuthmayer C, Neubauer S, Mairinger T, Zanghellini J and Hann S. ICT: isotope correction
20 1124 toolbox. *Bioinformatics*. 2016;32 1:154-6.
- 21 1125 77. Chokkathukalam A, Jankevics A, Creek DJ, Achcar F, Barrett MP and Breitling R. mzMatch-
22 1126 ISO: an R tool for the annotation and relative quantification of isotope-labelled mass
23 1127 spectrometry data. *Bioinformatics*. 2013;29 2:281-3.
- 24 1128 78. Bueschl C, Kluger B, Berthiller F, Lirk G, Winkler S, Krska R, et al. MetExtract: a new software
25 1129 tool for the automated comprehensive extraction of metabolite-derived LC/MS signals in
26 1130 metabolomics research. *Bioinformatics*. 2012;28 5:736-8.
- 27 1131 79. Huang X, Chen Y-J, Cho K, Nikolskiy I, Crawford PA and Patti GJ. X13CMS: global tracking of
28 1132 isotopic labels in untargeted metabolomics. *Analytical Chemistry*. 2014;86 3:1632-9.
- 29 1133 80. Capellades J, Navarro M, Samino S, Garcia-Ramirez M, Hernandez C, Simo R, et al. geoRge: A
30 1134 computational tool to detect the presence of stable isotope labeling in LC/MS-based
31 1135 untargeted metabolomics. *Analytical Chemistry*. 2015;88 1:621-8.
- 32 1136 81. Weindl D, Wegner A and Hiller K. MIA: non-targeted mass isotopologue analysis.
33 1137 *Bioinformatics*. 2016:btw317.
- 34 1138 82. Cai Y, Weng K, Guo Y, Peng J and Zhu Z-J. An integrated targeted metabolomic platform for
35 1139 high-throughput metabolite profiling and automated data processing. *Metabolomics*.
36 1140 2015;11 6:1575-86.
- 37 1141 83. Wong JW, Abuhusain HJ, McDonald KL and Don AS. MMSAT: automated quantification of
38 1142 metabolites in selected reaction monitoring experiments. *Analytical Chemistry*. 2011;84
39 1143 1:470-4.
- 40 1144 84. Tsugawa H, Arita M, Kanazawa M, Ogiwara A, Bamba T and Fukusaki E. MRMPROBS: A data
41 1145 assessment and metabolite identification tool for large-scale multiple reaction monitoring
42 1146 based widely targeted metabolomics. *Analytical Chemistry*. 2013;85 10:5191-9.
- 43 1147 85. Nikolskiy I, Mahieu NG, Chen Y-J, Tautenhahn R and Patti GJ. An untargeted metabolomic
44 1148 workflow to improve structural characterization of metabolites. *Analytical Chemistry*.
45 1149 2013;85 16:7713-9.
- 46 1150 86. Ma Y, Kind T, Yang D, Leon C and Fiehn O. MS2Analyzer: A software for small molecule
47 1151 substructure annotations from accurate tandem mass spectra. *Analytical Chemistry*. 2014;86
48 1152 21:10724-31.

58
59
60
61
62
63
64
65

- 1153 87. van der Hoof J, Wandy J, Barrett MP, Burgess KEV and Rogers S. Topic modeling for
1 1154 untargeted substructure exploration in metabolomics. *Proceedings of the National Academy
2 1155 of Sciences*. 2016;113 48:13738-43. doi:10.1073/pnas.1608041113.
- 3 1156 88. Tsugawa H, Cajka T, Kind T, Ma Y, Higgins B, Ikeda K, et al. MS-DIAL: data-independent
4 1157 MS/MS deconvolution for comprehensive metabolome analysis. *Nature methods*. 2015;12
5 1158 6:523-6.
- 7 1159 89. Li H, Cai Y, Guo Y, Chen F and Zhu Z-J. MetDIA: Targeted Metabolite Extraction of
8 1160 Multiplexed MS/MS Spectra Generated by Data-Independent Acquisition. *Analytical
9 1161 Chemistry*. 2016;88 17:8757-64.
- 10 1162 90. Libiseller G, Dvorzak M, Kleb U, Gander E, Eisenberg T, Madeo F, et al. IPO: a tool for
11 1163 automated optimization of XCMS parameters. *BMC Bioinformatics*. 2015;16 1:118.
- 13 1164 91. Mahieu NG, Huang X, Chen Y-J and Patti GJ. Credentialing features: a platform to benchmark
14 1165 and optimize untargeted metabolomic methods. *Analytical Chemistry*. 2014;86 19:9583-9.
- 15 1166 92. Brodsky L, Moussaieff A, Shahaf N, Aharoni A and Rogachev I. Evaluation of Peak Picking
16 1167 Quality in LC- MS Metabolomics Data. *Analytical Chemistry*. 2010;82 22:9177-87.
- 17 1168 93. Ranjbar MRN, Di Poto C, Wang Y and Ressom HW. Simat: Gc-sim-ms data analysis tool. *BMC
18 1169 Bioinformatics*. 2015;16 1:259.
- 20 1170 94. Mak TD, Laiakis EC, Goudarzi M and Fornace Jr AJ. Metabolizer: A novel statistical workflow
21 1171 for analyzing postprocessed lc-ms metabolomics data. *Analytical Chemistry*. 2013;86 1:506-
22 1172 13.
- 23 1173 95. Kastenmüller G, Römisch-Margl W, Wägele B, Altmaier E and Suhre K. metaP-server: a web-
24 1174 based metabolomics data analysis tool. *BioMed Research International*. 2010;2011.
- 26 1175 96. Fitzpatrick MA, McGrath CM and Young SP. Pathomx: an interactive workflow-based tool for
27 1176 the analysis of metabolomic data. *BMC Bioinformatics*. 2014;15 1:396.
- 28 1177 97. Hughes G, Cruickshank-Quinn C, Reisdorph R, Lutz S, Petrache I, Reisdorph N, et al.
29 1178 MSPrep—Summarization, normalization and diagnostics for processing of mass
30 1179 spectrometry-based metabolomic data. *Bioinformatics*. 2014;30 1:133-4.
- 32 1180 98. Sun X and Weckwerth W. COVAIn: a toolbox for uni-and multivariate statistics, time-series
33 1181 and correlation network analysis and inverse estimation of the differential Jacobian from
34 1182 metabolomics covariance data. *Metabolomics*. 2012;8 1:81-93.
- 35 1183 99. Glaab E and Schneider R. RepExplore: Addressing technical replicate variance in proteomics
36 1184 and metabolomics data analysis. *Bioinformatics*. 2015:btv127.
- 38 1185 100. Zhan X, Patterson AD and Ghosh D. Kernel approaches for differential expression analysis of
39 1186 mass spectrometry-based metabolomics data. *BMC Bioinformatics*. 2015;16 1:77.
- 40 1187 101. Nodzinski M, Muehlbauer MJ, Bain JR, Reisetter AC, Lowe WL and Scholtens DM.
41 1188 Metabomxtr: an R package for mixture-model analysis of non-targeted metabolomics data.
42 1189 *Bioinformatics*. 2014;30 22:3287-8.
- 44 1190 102. Suvitaival T, Rogers S and Kaski S. Stronger findings from mass spectral data through multi-
45 1191 peak modeling. *BMC Bioinformatics*. 2014;15 1:208.
- 46 1192 103. Mak TD, Laiakis EC, Goudarzi M and Fornace Jr AJ. Selective paired ion contrast analysis: a
47 1193 novel algorithm for analyzing postprocessed LC-MS metabolomics data possessing high
48 1194 experimental noise. *Analytical Chemistry*. 2015;87 6:3177-86.
- 50 1195 104. Ernest B, Gooding JR, Campagna SR, Saxton AM and Voy BH. MetabR: an R script for linear
51 1196 model analysis of quantitative metabolomic data. *BMC research notes*. 2012;5 1:596.
- 52 1197 105. Huang J-H, Yan J, Wu Q-H, Ferro MD, Yi L-Z, Lu H-M, et al. Selective of informative
53 1198 metabolites using random forests based on model population analysis. *Talanta*.
54 1199 2013;117:549-55.
- 56 1200 106. Simader AM, Kluger B, Neumann NKN, Bueschl C, Lemmens M, Lirk G, et al. QCScreen: a
57 1201 software tool for data quality control in LC-HRMS based metabolomics. *BMC Bioinformatics*.
58 1202 2015;16 1:341.

59
60
61
62
63
64
65

- 1203 107. Styczynski MP, Moxley JF, Tong LV, Walther JL, Jensen KL and Stephanopoulos GN.
1 1204 Systematic identification of conserved metabolites in GC/MS data for metabolomics and
2 1205 biomarker discovery. *Analytical Chemistry*. 2007;79 3:966-73.
- 3 1206 108. Baran R, Kochi H, Saito N, Suematsu M, Soga T, Nishioka T, et al. MathDAMP: a package for
4 1207 differential analysis of metabolite profiles. *BMC Bioinformatics*. 2006;7 1:530.
- 5 1208 109. Fernie AR. The future of metabolic phytochemistry: Larger numbers of metabolites, higher
6 1209 resolution, greater understanding. *Phytochemistry*. 2007;68 22–24:2861-80.
7 1210 doi:<http://dx.doi.org/10.1016/j.phytochem.2007.07.010>.
- 8 1211 110. Tohge T, Wendenburg R, Ishihara H, Nakabayashi R, Watanabe M, Sulpice R, et al.
9 1212 Characterization of a recently evolved flavonol-phenylacetyltransferase gene provides
10 1213 signatures of natural light selection in Brassicaceae. *Nature communications*. 2016;7.
- 11 1214 111. Zhou B, Wang J and Ransom HW. MetaboSearch: tool for mass-based metabolite
12 1215 identification using multiple databases. *PLoS One*. 2012;7 6:e40096.
- 13 1216 112. Brown M, Wedge DC, Goodacre R, Kell DB, Baker PN, Kenny LC, et al. Automated workflows
14 1217 for accurate mass-based putative metabolite identification in LC/MS-derived metabolomic
15 1218 datasets. *Bioinformatics*. 2011;27 8:1108-12.
- 16 1219 113. Daly R, Rogers S, Wandy J, Jankevics A, Burgess KE and Breitling R. MetAssign: probabilistic
17 1220 annotation of metabolites from LC–MS data using a Bayesian clustering approach.
18 1221 *Bioinformatics*. 2014;30 19:2764-71.
- 19 1222 114. Böcker S, Letzel MC, Lipták Z and Pervukhin A. SIRIUS: decomposing isotope patterns for
20 1223 metabolite identification. *Bioinformatics*. 2009;25 2:218-24.
- 21 1224 115. Sakurai N, Ara T, Kanaya S, Nakamura Y, Iijima Y, Enomoto M, et al. An application of a
22 1225 relational database system for high-throughput prediction of elemental compositions from
23 1226 accurate mass values. *Bioinformatics*. 2013;29 2:290-1.
- 24 1227 116. Lommen A. Ultrafast PubChem searching combined with improved filtering rules for
25 1228 elemental composition analysis. *Analytical Chemistry*. 2014;86 11:5463-9.
- 26 1229 117. Dhanasekaran AR, Pearson JL, Ganesan B and Weimer BC. Metabolome searcher: a high
27 1230 throughput tool for metabolite identification and metabolic pathway mapping directly from
28 1231 mass spectrometry and using genome restriction. *BMC Bioinformatics*. 2015;16 1:62.
- 29 1232 118. Suhre K and Schmitt-Kopplin P. MassTRIX: mass translator into pathways. *Nucleic acids*
30 1233 *research*. 2008;36 suppl 2:W481-W4.
- 31 1234 119. Uppal K, Soltow QA, Promislow DE, Wachtman LM, Quyyumi AA and Jones DP. MetabNet: an
32 1235 R package for metabolic association analysis of high-resolution metabolomics data. *Frontiers*
33 1236 *in bioengineering and biotechnology*. 2015;3:87.
- 34 1237 120. Silva RR, Jourdan F, Salvanha DM, Letisse F, Jamin EL, Guidetti-Gonzalez S, et al. ProbMetab:
35 1238 an R package for Bayesian probabilistic annotation of LC–MS-based metabolomics.
36 1239 *Bioinformatics*. 2014;30 9:1336-7.
- 37 1240 121. Rogers S, Scheltema RA, Girolami M and Breitling R. Probabilistic assignment of formulas to
38 1241 mass peaks in metabolomics experiments. *Bioinformatics*. 2009;25 4:512-8.
39 1242 doi:10.1093/bioinformatics/btn642.
- 40 1243 122. Weber RJ and Viant MR. MI-Pack: Increased confidence of metabolite identification in mass
41 1244 spectra by integrating accurate masses and metabolic pathways. *Chemometrics and*
42 1245 *Intelligent Laboratory Systems*. 2010;104 1:75-82.
- 43 1246 123. Qiu F, Fine DD, Wherritt DJ, Lei Z and Sumner LW. PlantMAT: A Metabolomics Tool for
44 1247 Predicting the Specialized Metabolic Potential of a System and for Large-Scale Metabolite
45 1248 Identifications. *Analytical Chemistry*. 2016;88 23:11373-83.
- 46 1249 124. Ruttkies C, Schymanski EL, Wolf S, Hollender J and Neumann S. MetFrag relaunched:
47 1250 incorporating strategies beyond in silico fragmentation. *Journal of cheminformatics*. 2016;8
48 1251 1:3.

58
59
60
61
62
63
64
65

- 1252 125. Menikarachchi LC, Cawley S, Hill DW, Hall LM, Hall L, Lai S, et al. MolFind: a software package enabling HPLC/MS-based identification of unknown chemical structures. *Analytical Chemistry*. 2012;84 21:9388-94.
- 1 1253
- 2 1254
- 3 1255 126. Allen F, Pon A, Wilson M, Greiner R and Wishart D. CFM-ID: a web server for annotation, spectrum prediction and metabolite identification from tandem mass spectra. *Nucleic acids research*. 2014;42 W1:W94-W9.
- 4 1256
- 5 1257
- 6 1258 127. Ridder L, van der Hoof JJ and Verhoeven S. Automatic compound annotation from mass spectrometry data using MAGMa. *Mass Spectrometry*. 2014;3 Special_Issue_2:S0033-S.
- 7 1259
- 8 1260 128. Heinonen M, Shen H, Zamboni N and Rousu J. Metabolite identification and molecular fingerprint prediction through machine learning. *Bioinformatics*. 2012;28 18:2333-41.
- 9 1261
- 10 1262 129. Dührkop K, Shen H, Meusel M, Rousu J and Böcker S. Searching molecular structure databases with tandem mass spectra using CSI: FingerID. *Proceedings of the National Academy of Sciences*. 2015;112 41:12580-5.
- 11 1263
- 12 1264 130. Gerlich M and Neumann S. MetFusion: integration of compound identification strategies. *Journal of Mass Spectrometry*. 2013;48 3:291-8.
- 13 1265
- 14 1266 131. Leader DP, Burgess K, Creek D and Barrett MP. Pathos: A web facility that uses metabolic maps to display experimental changes in metabolites identified by mass spectrometry. *Rapid Communications in Mass Spectrometry*. 2011;25 22:3422-6.
- 15 1267
- 16 1268 132. Pon A, Jewison T, Su Y, Liang Y, Knox C, Maciejewski A, et al. Pathways with PathWhiz. *Nucleic acids research*. 2015:gkv399.
- 17 1269
- 18 1270 133. Yamada T, Letunic I, Okuda S, Kanehisa M and Bork P. iPath2. 0: interactive pathway explorer. *Nucleic acids research*. 2011;39 suppl 2:W412-W5.
- 19 1271
- 20 1272 134. Kutmon M, van Iersel MP, Bohler A, Kelder T, Nunes N, Pico AR, et al. PathVisio 3: an extendable pathway analysis toolbox. *PLoS Comput Biol*. 2015;11 2:e1004085.
- 21 1273
- 22 1274 135. Pathan M, Keerthikumar S, Ang CS, Gangoda L, Quek CY, Williamson NA, et al. FunRich: An open access standalone functional enrichment and interaction network analysis tool. *Proteomics*. 2015;15 15:2597-601.
- 23 1275
- 24 1276 136. Moreno P, Beisken S, Harsha B, Muthukrishnan V, Tudose I, Dekker A, et al. BiNChE: a web tool and library for chemical enrichment analysis based on the ChEBI ontology. *BMC Bioinformatics*. 2015;16 1:56.
- 25 1277
- 26 1278 137. Kankainen M, Gopalacharyulu P, Holm L and Orešič M. MPEA—metabolite pathway enrichment analysis. *Bioinformatics*. 2011;27 13:1878-9.
- 27 1279
- 28 1280 138. Aggio RB, Ruggiero K and Villas-Bôas SG. Pathway Activity Profiling (PAPi): from the metabolite profile to the metabolic pathway activity. *Bioinformatics*. 2010;26 23:2969-76.
- 29 1281
- 30 1282 139. Eichner J, Rosenbaum L, Wrzodek C, Häring H-U, Zell A and Lehmann R. Integrated enrichment analysis and pathway-centered visualization of metabolomics, proteomics, transcriptomics, and genomics data by using the InCroMAP software. *Journal of Chromatography B*. 2014;966:77-82.
- 31 1283
- 32 1284 140. Carazzolle MF, de Carvalho LM, Slepicka HH, Vidal RO, Pereira GAG, Kobarg J, et al. IIS—Integrated Interactome System: a web-based platform for the annotation, analysis and visualization of protein-metabolite-gene-drug interactions by integrating a variety of data sources and tools. *PLoS One*. 2014;9 6:e100385.
- 33 1285
- 34 1286 141. Sakurai N, Ara T, Ogata Y, Sano R, Ohno T, Sugiyama K, et al. KaPPA-View4: a metabolic pathway database for representation and analysis of correlation networks of gene co-expression and metabolite co-accumulation and omics data. *Nucleic acids research*. 2011;39 suppl 1:D677-D84.
- 35 1287
- 36 1288 142. Usadel B, Poree F, Nagel A, Lohse M, CZEDIK-EYSENBERG A and Stitt M. A guide to using MapMan to visualize and compare Omics data in plants: a case study in the crop species, Maize. *Plant, cell & environment*. 2009;32 9:1211-29.
- 37 1289
- 38 1290 143. Neuweger H, Persicke M, Albaum SP, Bekel T, Dondrup M, Hüser AT, et al. Visualizing post genomics data-sets on customized pathway maps by ProMeTra—aeration-dependent gene
- 39 1291
- 40 1292
- 41 1293
- 42 1294
- 43 1295
- 44 1296
- 45 1297
- 46 1298
- 47 1299
- 48 1300
- 49 1301
- 50 1302
- 51
- 52
- 53
- 54
- 55
- 56
- 57
- 58
- 59
- 60
- 61
- 62
- 63
- 64
- 65

- 1303 expression and metabolism of *Corynebacterium glutamicum* as an example. *BMC systems biology*. 2009;3 1:82.
- 1 1304
- 2 1305 144. García-Alcalde F, García-López F, Dopazo J and Conesa A. Paintomics: a web based tool for
- 3 1306 the joint visualization of transcriptomics and metabolomics data. *Bioinformatics*. 2011;27
- 4 1307 1:137-9.
- 5 1308 145. Rohn H, Junker A, Hartmann A, Grafahrend-Belau E, Treutler H, Klapperstück M, et al.
- 6 1309 VANTED v2: a framework for systems biology applications. *BMC systems biology*. 2012;6
- 7 1310 1:139.
- 8 1311 146. López-Ibáñez J, Pazos F and Chagoyen M. MBROLE 2.0—functional enrichment of chemical
- 9 1312 compounds. *Nucleic acids research*. 2016;44 W1:W201-W4.
- 10 1313 147. Kamburov A, Cavill R, Ebbels TM, Herwig R and Keun HC. Integrated pathway-level analysis
- 11 1314 of transcriptomics and metabolomics data with IMPaLA. *Bioinformatics*. 2011;27 20:2917-8.
- 12 1315 148. Jourdan F, Breitling R, Barrett MP and Gilbert D. MetaNetter: inference and visualization of
- 13 1316 high-resolution metabolomic networks. *Bioinformatics*. 2008;24 1:143-5.
- 14 1317 149. Grapov D, Wanichthanarak K and Fiehn O. MetaMapR: pathway independent metabolomic
- 15 1318 network analysis incorporating unknowns. *Bioinformatics*. 2015:btv194.
- 16 1319 150. Lu J and Carlson HA. ChemTreeMap: an interactive map of biochemical similarity in
- 17 1320 molecular datasets. *Bioinformatics*. 2016;32 23:3584-92.
- 18 1321 doi:10.1093/bioinformatics/btw523.
- 19 1322 151. Treutler H, Tsugawa H, Porzel A, Gorzolka K, Tissier A, Neumann S, et al. Discovering
- 20 1323 Regulated Metabolite Families in Untargeted Metabolomics Studies. *Analytical Chemistry*.
- 21 1324 2016;88 16:8082-90.
- 22 1325 152. Hamdalla MA, Rajasekaran S, Grant DF and Măndoiu II. Metabolic Pathway Predictions for
- 23 1326 Metabolomics: A Molecular Structure Matching Approach. *Journal of chemical information*
- 24 1327 *and modeling*. 2015;55 3:709-18.
- 25 1328 153. Pence HE and Williams A. ChemSpider: an online chemical information resource. ACS
- 26 1329 Publications, 2010.
- 27 1330 154. Kim S, Thiessen PA, Bolton EE, Chen J, Fu G, Gindulyte A, et al. PubChem substance and
- 28 1331 compound databases. *Nucleic acids research*. 2015:gkv951.
- 29 1332 155. Hastings J, Owen G, Dekker A, Ennis M, Kale N, Muthukrishnan V, et al. ChEBI in 2016:
- 30 1333 Improved services and an expanding collection of metabolites. *Nucleic acids research*.
- 31 1334 2015:gkv1031.
- 32 1335 156. Gaulton A, Bellis LJ, Bento AP, Chambers J, Davies M, Hersey A, et al. ChEMBL: a large-scale
- 33 1336 bioactivity database for drug discovery. *Nucleic acids research*. 2012;40 D1:D1100-D7.
- 34 1337 157. Seiler KP, George GA, Happ MP, Bodycombe NE, Carrinski HA, Norton S, et al. ChemBank: a
- 35 1338 small-molecule screening and cheminformatics resource database. *Nucleic acids research*.
- 36 1339 2008;36 suppl 1:D351-D9.
- 37 1340 158. Wishart DS, Jewison T, Guo AC, Wilson M, Knox C, Liu Y, et al. HMDB 3.0—the human
- 38 1341 metabolome database in 2013. *Nucleic acids research*. 2012:gks1065.
- 39 1342 159. Cui Q, Lewis IA, Hegeman AD, Anderson ME, Li J, Schulte CF, et al. Metabolite identification
- 40 1343 via the Madison Metabolomics Consortium Database. *Nat Biotech*. 2008;26 2:162-4.
- 41 1344 doi:http://www.nature.com/nbt/journal/v26/n2/supinfo/nbt0208-162_S1.html.
- 42 1345 160. Masciocchi J, Frau G, Fanton M, Sturlese M, Floris M, Pireddu L, et al. MMsINC: a large-scale
- 43 1346 cheminformatics database. *Nucleic acids research*. 2009;37 suppl 1:D284-D90.
- 44 1347 161. Afendi FM, Okada T, Yamazaki M, Hirai-Morita A, Nakamura Y, Nakamura K, et al. KNApSACK
- 45 1348 family databases: integrated metabolite–plant species databases for multifaceted plant
- 46 1349 research. *Plant and Cell Physiology*. 2012;53 2:e1-e.
- 47 1350 162. Arita M and Suwa K. Search extension transforms Wiki into a relational system: a case for
- 48 1351 flavonoid metabolite database. *BioData mining*. 2008;1 1:7.
- 49
- 50
- 51
- 52
- 53
- 54
- 55
- 56
- 57
- 58
- 59
- 60
- 61
- 62
- 63
- 64
- 65

- 1352 163. Sharma A, Dutta P, Sharma M, Rajput NK, Dodiya B, George JJ, et al. BioPhytMol: a drug
1 1353 discovery community resource on anti-mycobacterial phytomolecules and plant extracts.
2 1354 Journal of cheminformatics. 2014;6 1:46.
- 3 1355 164. Kumari S, Pundhir S, Priya P, Jeena G, Punetha A, Chawla K, et al. EssOilDB: a database of
4 1356 essential oils reflecting terpene composition and variability in the plant kingdom. Database.
5 1357 2014;2014:bau120.
- 6 1357 165. Hummel J, Selbig J, Walther D and Kopka J. The Golm Metabolome Database: a database for
7 1358 GC-MS based metabolite profiling. Metabolomics. Springer; 2007. p. 75-95.
- 8 1359 166. Skogerson K, Wohlgemuth G, Barupal DK and Fiehn O. The volatile compound BinBase mass
9 1360 spectral database. BMC Bioinformatics. 2011;12 1:321.
- 10 1361 167. Horai H, Arita M, Kanaya S, Nihei Y, Ikeda T, Suwa K, et al. MassBank: a public repository for
11 1362 sharing mass spectral data for life sciences. Journal of Mass Spectrometry. 2010;45 7:703-14.
- 12 1363 168. Smith CA, O'Maille G, Want EJ, Qin C, Trauger SA, Brandon TR, et al. METLIN: a metabolite
13 1364 mass spectral database. Therapeutic drug monitoring. 2005;27 6:747-51.
- 14 1365 169. Cho K, Mahieu N, Ivanisevic J, Uritboonthai W, Chen Y-J, Siuzdak G, et al. isoMETLIN: a
15 1366 database for isotope-based metabolomics. Analytical Chemistry. 2014;86 19:9358-61.
- 16 1367 170. Wishart D, Arndt D, Pon A, Sajed T, Guo AC, Djoumbou Y, et al. T3DB: the toxic exposome
17 1368 database. Nucleic acids research. 2015;43 D1:D928-D34.
- 18 1369 171. Cuthbertson DJ, Johnson SR, Piljac-Žegarac J, Kappel J, Schäfer S, Wüst M, et al. Accurate
19 1370 mass-time tag library for LC/MS-based metabolite profiling of medicinal plants.
20 1371 Phytochemistry. 2013;91:187-97.
- 21 1372 172. Sawada Y, Nakabayashi R, Yamada Y, Suzuki M, Sato M, Sakata A, et al. RIKEN tandem mass
22 1373 spectral database (ReSpect) for phytochemicals: a plant-specific MS/MS-based data resource
23 1374 and database. Phytochemistry. 2012;82:38-45.
- 24 1375 173. Shahaf N, Rogachev I, Heinig U, Meir S, Malitsky S, Battat M, et al. The WEIZMASS spectral
25 1376 library for high-confidence metabolite identification. Nature communications. 2016;7.
- 26 1377 174. Weber RJM, Li E, Bruty J, He S and Viant MR. MaConDa: a publicly accessible mass
27 1378 spectrometry contaminants database. Bioinformatics. 2012;28 21:2856-7.
28 1379 doi:10.1093/bioinformatics/bts527.
- 29 1380 175. Matsuda F, Hirai MY, Sasaki E, Akiyama K, Yonekura-Sakakibara K, Provart NJ, et al.
30 1381 AtMetExpress development: a phytochemical atlas of Arabidopsis development. Plant
31 1382 Physiology. 2010;152 2:566-78.
- 32 1383 176. Fukushima A, Kusano M, Mejia RF, Iwasa M, Kobayashi M, Hayashi N, et al. Metabolomic
33 1384 characterization of knockout mutants in Arabidopsis: development of a metabolite profiling
34 1385 database for knockout mutants in Arabidopsis. Plant Physiology. 2014;165 3:948-61.
- 35 1386 177. Moco S, Bino RJ, Vorst O, Verhoeven HA, de Groot J, van Beek TA, et al. A liquid
36 1387 chromatography-mass spectrometry-based metabolome database for tomato. Plant
37 1388 Physiology. 2006;141 4:1205-18.
- 38 1389 178. Brockmöller T, Ling Z, Li D, Gaquerel E, Baldwin IT and Xu S. Nicotiana attenuata Data Hub
39 1390 (Na DH): an integrative platform for exploring genomic, transcriptomic and metabolomic
40 1391 data in wild tobacco. BMC genomics. 2017;18 1:79.
- 41 1392 179. Colmsee C, Mascher M, Czauderna T, Hartmann A, Schlüter U, Zellerhoff N, et al. OPTIMAS-
42 1393 DW: a comprehensive transcriptomics, metabolomics, ionomics, proteomics and phenomics
43 1394 data resource for maize. BMC plant biology. 2012;12 1:245.
- 44 1395 180. Joshi T, Yao Q, Levi DF, Brechenmacher L, Valliyodan B, Stacey G, et al. SoyMetDB: the
45 1396 soybean metabolome database. In: *Bioinformatics and Biomedicine (BIBM), 2010 IEEE
46 1397 International Conference on 2010*, pp.203-8. IEEE.
- 47 1398 181. Iijima Y, Nakamura Y, Ogata Y, Tanaka Ki, Sakurai N, Suda K, et al. Metabolite annotations
48 1399 based on the integration of mass spectral information. The Plant Journal. 2008;54 5:949-62.
- 49 1400

50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1401 182. Matsuda F, Yonekura-Sakakibara K, Niida R, Kuromori T, Shinozaki K and Saito K. MS/MS
1 1402 spectral tag-based annotation of non-targeted profile of plant secondary metabolites. *The*
2 1403 *Plant Journal*. 2009;57 3:555-77.
- 3 1404 183. Hur M, Campbell AA, Almeida-de-Macedo M, Li L, Ransom N, Jose A, et al. A global approach
4 1405 to analysis and interpretation of metabolic data for plant natural product discovery. *Natural*
5 1406 *product reports*. 2013;30 4:565-83.
- 7 1407 184. Wurtele ES, Chappell J, Jones AD, Celiz MD, Ransom N, Hur M, et al. Medicinal plants: a
8 1408 public resource for metabolomics and hypothesis development. *Metabolites*. 2012;2 4:1031-
9 1409 59.
- 10 1410 185. Sud M, Fahy E, Cotter D, Azam K, Vadivelu I, Burant C, et al. Metabolomics Workbench: An
11 1411 international repository for metabolomics data and metadata, metabolite standards,
12 1412 protocols, tutorials and training, and analysis tools. *Nucleic acids research*. 2015:gkv1042.
- 14 1413 186. Haug K, Salek RM, Conesa P, Hastings J, de Matos P, Rijnbeek M, et al. MetaboLights—an
15 1414 open-access general-purpose repository for metabolomics studies and associated meta-
16 1415 data. *Nucleic acids research*. 2012:gks1004.
- 17 1416 187. Cook CE, Bergman MT, Finn RD, Cochrane G, Birney E and Apweiler R. The European
18 1417 Bioinformatics Institute in 2016: data growth and integration. *Nucleic acids research*.
20 1418 2016;44 D1:D20-D6.
- 21 1419 188. Wang M, Carver JJ, Phelan VV, Sanchez LM, Garg N, Peng Y, et al. Sharing and community
22 1420 curation of mass spectrometry data with Global Natural Products Social Molecular
23 1421 Networking. *Nat Biotechnol*. 2016;34 8:828-37.
- 24 1422 189. Kanehisa M and Goto S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic acids*
25 1423 *research*. 2000;28 1:27-30.
- 27 1424 190. Kelder T, Pico AR, Hanspers K, Van Iersel MP, Evelo C and Conklin BR. Mining biological
28 1425 pathways using WikiPathways web services. *PLoS One*. 2009;4 7:e6447.
- 29 1426 191. Navas-Delgado I, García-Godoy MJ, López-Camacho E, Rybinski M, Reyes-Palomares A,
30 1427 Medina MÁ, et al. kpath: integration of metabolic pathway linked data. *Database*.
32 1428 2015;2015:bav053.
- 33 1429 192. Caspi R, Foerster H, Fulcher CA, Kaipa P, Krummenacker M, Latendresse M, et al. The
34 1430 MetaCyc Database of metabolic pathways and enzymes and the BioCyc collection of
35 1431 Pathway/Genome Databases. *Nucleic acids research*. 2008;36 suppl 1:D623-D31.
- 36 1432 193. Arkin AP, Stevens RL, Cottingham RW, Maslov S, Henry CS, Dehal P, et al. The DOE Systems
37 1433 Biology Knowledgebase (KBBase). *bioRxiv*. 2016:096354.
- 39 1434 194. Schreiber F, Colmsee C, Czauderna T, Grafahrend-Belau E, Hartmann A, Junker A, et al.
40 1435 MetaCrop 2.0: managing and exploring information about crop plant metabolism. *Nucleic*
41 1436 *acids research*. 2011:gkr1004.
- 42 1437 195. Sucaet Y, Wang Y, Li J and Wurtele ES. MetNet Online: a novel integrated resource for plant
43 1438 systems biology. *BMC Bioinformatics*. 2012;13 1:267.
- 45 1439 196. Tello-Ruiz MK, Stein J, Wei S, Preece J, Olson A, Naithani S, et al. Gramene 2016:
46 1440 comparative plant genomics and pathway resources. *Nucleic acids research*. 2015:gkv1179.
- 47 1441 197. Ott MA and Vriend G. Correcting ligands, metabolites, and pathways. *BMC Bioinformatics*.
48 1442 2006;7 1:517.
- 49 1443 198. Lang M, Stelzer M and Schomburg D. BKM-react, an integrated biochemical reaction
50 1444 database. *BMC biochemistry*. 2011;12 1:42.
- 52 1445 199. Kumar A, Suthers PF and Maranas CD. MetRxn: a knowledgebase of metabolites and
53 1446 reactions spanning metabolic models and databases. *BMC Bioinformatics*. 2012;13 1:6.
- 54 1447 200. Jeffryes JG, Colastani RL, Elbadawi-Sidhu M, Kind T, Niehaus TD, Broadbelt LJ, et al. MINEs:
55 1448 open access databases of computationally predicted enzyme promiscuity products for
56 1449 untargeted metabolomics. *Journal of cheminformatics*. 2015;7 1:44.
- 58 1450 201. Zhou Z, Shen X, Tu J and Zhu Z-J. Large-Scale Prediction of Collision Cross-Section Values for
59 1451 Metabolites in Ion Mobility-Mass Spectrometry. *Analytical Chemistry*. 2016;88 22:11084-91.

60
61
62
63
64
65

- 1452 202. Allard P-M, Péresse T, Bisson J, Gindro K, Marcourt L, Pham VC, et al. Integration of
1 1453 molecular networking and in-silico MS/MS fragmentation for natural products dereplication.
2 1454 Analytical Chemistry. 2016;88 6:3317-23.
- 3 1455 203. Palla P, Frau G, Vargiu L and Rodriguez-Tomé P. QTREDS: a Ruby on Rails-based platform for
4 1456 omics laboratories. BMC Bioinformatics. 2014;15 1:S13.
- 5 1457 204. Hunter A, Dayalan S, De Souza D, Power B, Lorrimar R, Szabo T, et al. MASTR-MS: a web-
6 1458 based collaborative laboratory information management system (LIMS) for metabolomics.
7 1459 Metabolomics. 2017;13 2:14.
- 8 1459 205. Franceschi P, Mylonas R, Shahaf N, Scholz M, Arapitsas P, Masuero D, et al. MetaDB a data
9 1460 processing workflow in untargeted MS-based metabolomics experiments. Frontiers in
10 1461 bioengineering and biotechnology. 2014;2:72.
- 11 1462 206. Ara T, Enomoto M, Arita M, Ikeda C, Kera K, Yamada M, et al. Metabolonote: a wiki-based
12 1463 database for managing hierarchical metadata of metabolome analyses. Frontiers in
13 1464 bioengineering and biotechnology. 2015;3:38.
- 14 1465 207. Wohlgemuth G, Mehta SS, Mejia RF, Neumann S, Pedrosa D, Pluskal T, et al. SPLASH, a
15 1466 hashed identifier for mass spectra. Nat Biotechnol. 2016;34 11:1099-101.
- 16 1467 208. Redestig H, Kusano M, Fukushima A, Matsuda F, Saito K and Arita M. Consolidating
17 1468 metabolite identifiers to enable contextual and multi-platform metabolomics data analysis.
18 1469 BMC Bioinformatics. 2010;11 1:214.
- 19 1470 209. Wohlgemuth G, Haldiya PK, Willighagen E, Kind T and Fiehn O. The Chemical Translation
20 1471 Service—a web-based tool to improve standardization of metabolomic reports.
21 1472 Bioinformatics. 2010;26 20:2647-8.
- 22 1473 210. Carroll AJ, Zhang P, Whitehead L, Kaines S, Tcherkez G and Badger MR. PhenoMeter: a
23 1474 metabolome database search tool using statistical similarity matching of metabolic
24 1475 phenotypes for high-confidence detection of functional links. Frontiers in bioengineering and
25 1476 biotechnology. 2015;3.
- 26 1477 211. Sartor MA, Ade A, Wright Z, Omenn GS, Athey B and Karnovsky A. Metab2MeSH: annotating
27 1478 compounds with medical subject headings. Bioinformatics. 2012;28 10:1408-10.
- 28 1479 212. Xu QW, Griss J, Wang R, Jones AR, Hermjakob H and Vizcaíno JA. jmzTab: A Java interface to
29 1480 the mzTab data standard. Proteomics. 2014;14 11:1328-32.
- 30 1481 213. Scheltema RA, Jankevics A, Jansen RC, Swertz MA and Breitling R. PeakML/mzMatch: a file
31 1482 format, Java library, R library, and tool-chain for mass spectrometry data analysis. Analytical
32 1483 Chemistry. 2011;83 7:2786-93.
- 33 1484 214. Avtonomov DM, Raskind A and Nesvizhskii AI. BatMass: a Java Software Platform for LC–MS
34 1485 Data Visualization in Proteomics and Metabolomics. Journal of Proteome Research. 2016;15
35 1486 8:2500-9.
- 36 1487 215. Tanaka S, Fujita Y, Parry HE, Yoshizawa AC, Morimoto K, Murase M, et al. Mass++: A
37 1488 visualization and analysis tool for mass spectrometry. Journal of Proteome Research.
38 1489 2014;13 8:3846-53.
- 39 1490 216. Beisken S, Conesa P, Haug K, Salek RM and Steinbeck C. SpeckTackle: JavaScript charts for
40 1491 spectroscopy. Journal of cheminformatics. 2015;7 1:17.
- 41 1492 217. Stravs MA, Schymanski EL, Singer HP and Hollender J. Automatic recalibration and
42 1493 processing of tandem mass spectra using formula annotation. Journal of Mass Spectrometry.
43 1494 2013;48 1:89-99.
- 44 1495 218. Dong Y, Li B and Aharoni A. More than Pictures: When MS Imaging Meets Histology. Trends
45 1496 in plant science. 2016;21 8:686-98.
- 46 1497 219. Wijetunge CD, Saeed I, Boughton BA, Spraggins JM, Caprioli RM, Bacic A, et al. EXIMS: an
47 1498 improved data analysis pipeline based on a new peak picking method for EXploring Imaging
48 1499 Mass Spectrometry data. Bioinformatics. 2015;31 19:3198-206.
- 50 1500

58
59
60
61
62
63
64
65

- 1501 220. Rübél O, Greiner A, Cholia S, Louie K, Bethel EW, Northen TR, et al. OpenMSI: a high-
1 1502 performance web-based platform for mass spectrometry imaging. *Analytical Chemistry*.
2 1503 2013;85 21:10354-61.
- 3 1504 221. Husen P, Tarasov K, Katafiasz M, Sokol E, Vogt J, Baumgart J, et al. Analysis of lipid
4 1505 experiments (ALEX): a software framework for analysis of high-resolution shotgun lipidomics
5 1506 data. *PLoS One*. 2013;8 11:e79736.
- 7 1507 222. Tsugawa H, Ohta E, Izumi Y, Ogiwara A, Yukihiro D, Bamba T, et al. MRM-DIFF: data
8 1508 processing strategy for differential analysis in large scale MRM-based lipidomics studies.
9 1509 *Frontiers in genetics*. 2014;5.
- 10 1510 223. Wong G, Chan J, Kingwell BA, Leckie C and Meikle PJ. LICRE: unsupervised feature correlation
11 1511 reduction for lipidomics. *Bioinformatics*. 2014:btu381.
- 12 1512 224. Herzog R, Schuhmann K, Schwudke D, Sampaio JL, Bornstein SR, Schroeder M, et al.
13 1513 LipidXplorer: a software for consensual cross-platform lipidomics. *PLoS One*. 2012;7
14 1514 1:e29851.
- 16 1515 225. Haimi P, Uphoff A, Hermansson M and Somerharju P. Software tools for analysis of mass
17 1516 spectrometric lipidome data. *Analytical Chemistry*. 2006;78 24:8324-31.
- 19 1517 226. Blanchard AP, McDowell GS, Valenzuela N, Xu H, Gelbard S, Bertrand M, et al. Visualization
20 1518 and Phospholipid Identification (VaLID): online integrated search engine capable of
21 1519 identifying and visualizing glycerophospholipids with given mass. *Bioinformatics*. 2013;29
22 1520 2:284-5.
- 23 1521 227. Collins JR, Edwards BR, Fredricks HF and Van Mooy BA. LOBSTAHS: an adduct-based
24 1522 lipidomics strategy for discovery and identification of oxidative stress biomarkers. *Analytical
25 1523 Chemistry*. 2016;88 14:7154-62.
- 27 1524 228. Ahmed Z, Mayr M, Zeeshan S, Dandekar T, Mueller MJ and Fekete A. Lipid-Pro: a
28 1525 computational lipid identification solution for untargeted lipidomics on data-independent
29 1526 acquisition tandem mass spectrometry platforms. *Bioinformatics*. 2015;31 7:1150-3.
- 30 1527 229. Hartler J, Trötz Müller M, Chitruju C, Spener F, Köfeler HC and Thallinger GG. Lipid Data
31 1528 Analyzer: unattended identification and quantitation of lipids in LC-MS data. *Bioinformatics*.
32 1529 2011;27 4:572-7.
- 34 1530 230. Song H, Hsu F-F, Ladenson J and Turk J. Algorithm for processing raw mass spectrometric
35 1531 data to identify and quantitate complex lipid molecular species in mixtures by data-
36 1532 dependent scanning and fragment ion database searching. *Journal of the American Society
37 1533 for Mass Spectrometry*. 2007;18 10:1848-58.
- 39 1534 231. Sud M, Fahy E, Cotter D, Brown A, Dennis EA, Glass CK, et al. Lmsd: lipid maps structure
40 1535 database. *Nucleic acids research*. 2007;35 suppl 1:D527-D32.
- 41 1536 232. Watanabe K, Yasugi E and Oshima M. How to search the glycolipid data in "LIPIDBANK for
42 1537 Web", the newly developed lipid database in Japan. *Trends in Glycoscience and
43 1538 Glycotechnology*. 2000;12 65:175-84.
- 45 1539 233. Kind T, Liu K-H, Lee DY, DeFelice B, Meissen JK and Fiehn O. LipidBlast in silico tandem mass
46 1540 spectrometry database for lipid identification. *Nature methods*. 2013;10 8:755-8.
- 47 1541 234. Foster JM, Moreno P, Fabregat A, Hermjakob H, Steinbeck C, Apweiler R, et al. LipidHome: a
48 1542 database of theoretical lipids optimized for high throughput mass spectrometry lipidomics.
49 1543 *PLoS One*. 2013;8 5:e61951.
- 50 1544 235. Aimo L, Liechti R, Nospikel N, Niknejad A, Gleizes A, Götz L, et al. The SwissLipids
51 1545 knowledgebase for lipid biology. *Bioinformatics*. 2015:btv285.
- 53 1546 236. Tautenhahn R, Patti GJ, Rinehart D and Siuzdak G. XCMS Online: a web-based platform to
54 1547 process untargeted metabolomic data. *Analytical Chemistry*. 2012;84 11:5035-9.
- 55 1548 237. Grace SC, Embry S and Luo H. Haystack, a web-based tool for metabolomics research. *BMC
56 1549 Bioinformatics*. 2014;15 11:S12.
- 58 1550 238. Liang Y-J, Lin Y-T, Chen C-W, Lin C-W, Chao K-M, Pan W-H, et al. SMART: Statistical
59 1551 Metabolomics Analysis An R Tool. *Analytical Chemistry*. 2016;88 12:6334-41.
- 60
61
62
63
64
65

1552 239. Pluskal T, Castillo S, Villar-Briones A and Orešič M. MZmine 2: modular framework for
1 1553 processing, visualizing, and analyzing mass spectrometry-based molecular profile data. BMC
2 1554 Bioinformatics. 2010;11 1:395.

3 1555 240. Wei X, Sun W, Shi X, Koo I, Wang B, Zhang J, et al. MetSign: a computational platform for
4 1556 high-resolution mass spectrometry-based metabolomics. Analytical Chemistry. 2011;83
5 1557 20:7668-75.

7 1558 241. LaMarche BL, Crowell KL, Jaitly N, Petyuk VA, Shah AR, Polpitiya AD, et al. MultiAlign: a
8 1559 multiple LC-MS analysis tool for targeted omics analysis. BMC Bioinformatics. 2013;14 1:49.

9 1560 242. Carroll AJ, Badger MR and Millar AH. The MetabolomeExpress Project: enabling web-based
10 1561 processing, analysis and transparent dissemination of GC/MS metabolomics datasets. BMC
11 1562 Bioinformatics. 2010;11 1:376.

13 1563 243. Fernández-Albert F, Llorach R, Andrés-Lacueva C and Perera A. An R package to analyse
14 1564 LC/MS metabolomic data: MAIT (Metabolite Automatic Identification Toolkit).
15 1565 Bioinformatics. 2014;30 13:1937-9.

16 1566 244. Melamud E, Vastag L and Rabinowitz JD. Metabolomic analysis and visualization engine for
17 1567 LC- MS data. Analytical Chemistry. 2010;82 23:9818-26.

19 1568 245. Neuweger H, Albaum SP, Dondrup M, Persicke M, Watt T, Niehaus K, et al. MeltDB: a
20 1569 software platform for the analysis and integration of metabolomics experiment data.
21 1570 Bioinformatics. 2008;24 23:2726-32.

22 1571 246. Xia J, Sinelnikov IV, Han B and Wishart DS. MetaboAnalyst 3.0—making metabolomics more
23 1572 meaningful. Nucleic acids research. 2015;43 W1:W251-W7.

25 1573 247. Kaefer A, Landesfeind M, Feussner K, Mosblech A, Heilmann I, Morgenstern B, et al. MarVis-
26 1574 Pathway: integrative and exploratory pathway analysis of non-targeted metabolomics data.
27 1575 Metabolomics. 2015;11 3:764-77.

28 1576 248. Edmands WM, Barupal DK and Scalbert A. MetMSLine: an automated and fully integrated
29 1577 pipeline for rapid processing of high-resolution LC-MS metabolomic datasets. Bioinformatics.
30 1578 2014:btu705.

32 1579 249. Beisken S, Earll M, Portwood D, Seymour M and Steinbeck C. MassCascade: Visual
33 1580 Programming for LC-MS Data Processing in Metabolomics. Molecular informatics. 2014;33
34 1581 4:307-10.

35 1582 250. Winkler R. MASSyPup—an ‘Out of the Box’ solution for the analysis of mass spectrometry
36 1583 data. Journal of Mass Spectrometry. 2014;49 1:37-42.

38 1584 251. Sakurai N, Ara T, Enomoto M, Motegi T, Morishita Y, Kurabayashi A, et al. Tools and
39 1585 databases of the KOMICS web portal for preprocessing, mining, and dissemination of
40 1586 metabolomics data. BioMed Research International. 2014;2014.

41 1587 252. Sakurai T, Yamada Y, Sawada Y, Matsuda F, Akiyama K, Shinozaki K, et al. PRIME update:
42 1588 innovative content for plant metabolomics and integration of gene expression and
43 1589 metabolite accumulation. Plant and Cell Physiology. 2013;54 2:e5-e.

45 1590 253. Henry VJ, Bandrowski AE, Pepin A-S, Gonzalez BJ and Desfeux A. OMICtools: an informative
46 1591 directory for multi-omic data analysis. Database. 2014;2014:bau069.

47 1592 254. Sumner LW, Amberg A, Barrett D, Beale MH, Beger R, Daykin CA, et al. Proposed minimum
48 1593 reporting standards for chemical analysis. Metabolomics. 2007;3 3:211-21.
49 1594 doi:10.1007/s11306-007-0082-2.

51 1595 255. Gago J, Daloso DdM, Figueroa CM, Flexas J, Fernie AR and Nikoloski Z. Relationships of Leaf
52 1596 Net Photosynthesis, Stomatal Conductance, and Mesophyll Conductance to Primary
53 1597 Metabolism: A Multispecies Meta-Analysis Approach. Plant Physiology. 2016;171 1:265-79.
54 1598 doi:10.1104/pp.15.01660.

56 1599
57
58
59 1600 **Figure 1** Typical mass spectrometry based metabolomics workflow.

60
61
62
63
64
65

1601 **Additional file 1.xls** Summary of resources for mass spectrometry based metabolomics.

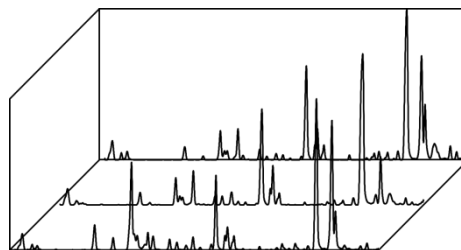
- 1
- 2
- 3
- 4
- 5
- 6
- 7
- 8
- 9
- 10
- 11
- 12
- 13
- 14
- 15
- 16
- 17
- 18
- 19
- 20
- 21
- 22
- 23
- 24
- 25
- 26
- 27
- 28
- 29
- 30
- 31
- 32
- 33
- 34
- 35
- 36
- 37
- 38
- 39
- 40
- 41
- 42
- 43
- 44
- 45
- 46
- 47
- 48
- 49
- 50
- 51
- 52
- 53
- 54
- 55
- 56
- 57
- 58
- 59
- 60
- 61
- 62
- 63
- 64
- 65

Sample Preparation

Data Acquisition

Processing

- Feature detection
- Alignment
- Quantification
- Spectra deconvolution
- Normalization

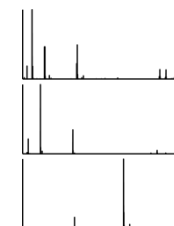


Features

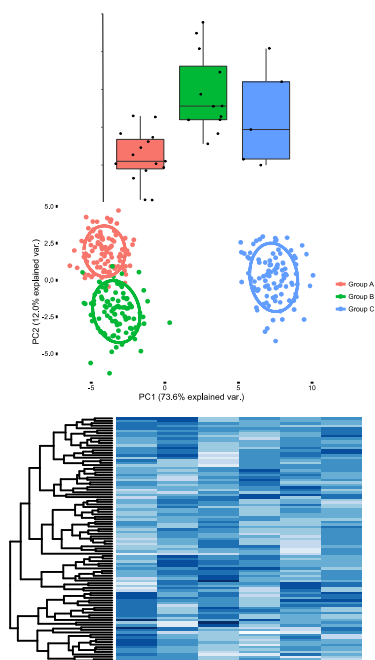
Samples

Intensities

Compound Spectra

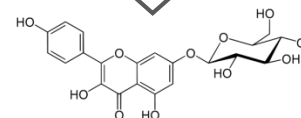
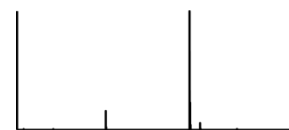


Statistical Analysis



Annotation

- Exact mass
- MSⁿ
- Spectra matching
- *In silico* prediction



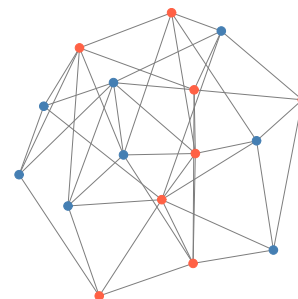
Databases

- Compounds
- Mass spectra
- Samples
- Pathway



Interpretation

- Network structure
- Pathway enrichment
- Integration





Click here to access/download
Supplementary Material
Databases_Submit.xls

