

# 1 **FragHub: A mass spectral libraries data integration workflow**

2  
3 Axel Dablan<sup>1,2,‡</sup>, Solweig Hennechart<sup>1,2,3,‡</sup>, Amélie Perez<sup>1,2</sup>, Guillaume Cabanac<sup>3,4</sup>, Yann  
4 Guitton<sup>5</sup>, Nils Paulhe<sup>6</sup>, Bernard Lyan<sup>6</sup>, Emilien Jamin<sup>7,2</sup>, Franck Giacomoni<sup>6</sup>, Guillaume  
5 Marti<sup>1,2\*</sup>

6  
7 <sup>1</sup>Laboratoire de Recherche en Sciences Végétales, Metatoul-AgromiX Platform, Université de Toulouse, CNRS,  
8 INP, 24 Chemin de Borde Rouge, Auzeville, 31320, Auzeville-Tolosane, France

9 <sup>2</sup>MetaboHUB-MetaToul, National Infrastructure of Metabolomics and Fluxomics, Toulouse, France

10 <sup>3</sup>Université Toulouse 3 – Paul Sabatier, IRIT UMR 5505 CNRS, Toulouse, France

11 <sup>4</sup>Institut Universitaire de France (IUF), Paris

12 <sup>5</sup>Oniris, INRAE, Laberca, 44300 Nantes, France

13 <sup>6</sup>Université Clermont Auvergne, INRAE, UNH, Plateforme d'Exploration du Métabolisme, MetaboHUB  
14 Clermont, Clermont-Ferrand F-63000, France

15 <sup>7</sup>Toxalim (Research Centre in Food Toxicology), Université de Toulouse, INRAE, ENVT, INP-Purpan, UPS,  
16 Toulouse, France

17 ‡These authors contributed equally

---

## 19 **Abstract**

20 Open mass spectral libraries (OMSL) are critical for metabolite annotation and machine learning, especially given the  
21 rising volume of untargeted metabolomic studies and the development of annotation pipelines. Despite their importance,  
22 the practical application of OMSLs is hampered by the lack of standardized file formats, metadata fields, and supporting  
23 ontology. Current libraries, often restricted to specific topics or matrices such as natural products, lipids, or the human  
24 metabolome, may limit the discovery potential of untargeted studies. FragHub addresses these challenges by integrating  
25 multiple OMSLs into a single comprehensive database, supporting various data formats and harmonizing metadata. It  
26 also proposes some generic filters for mass spectrum using a graphical user interface. Additionally, a workflow to  
27 generate in-house libraries compatible with FragHub is proposed. FragHub dynamically segregates libraries based on  
28 ionization modes and chromatography techniques, thereby enhancing data utility in metabolomic research. The FragHub  
29 Python code is publicly available under a MIT license, at the following repository:  
30 <https://github.com/eMetaboHUB/FragHub>. Generated data can be accessed at  
31 <https://doi.org/10.5281/zenodo.11057687>.

---

## 34 *Keywords*

35 *Open mass spectral library, metabolomics, dereplication, Mass spectrometry, database*

## 37 Introduction

38 Liquid Chromatography-Mass spectrometry (LC-MS) chemical profiling provides hundreds to  
39 thousands of features ( $m/z \times RT$  pairs) from a single biological matrix. The process of dereplication,  
40 which involves annotating all detected spectral signatures, is a major bottleneck in LC-MS based  
41 metabolomics<sup>1</sup>. Annotations rely on a “body of evidence” approach initially formalized by the  
42 Metabolomics Standards Initiative, stratified into four confidence levels: level 1, identified metabolites  
43 using authentic standard compounds; level 2, putatively annotated metabolites using public/commercial  
44 spectral libraries; level 3, putatively characterized metabolites based on diagnostic ions and/or partial  
45 spectral similarities to known compounds of a chemical class; and level 4, unknown metabolites<sup>2</sup>. These  
46 confidence levels have been further refined to include new strategies such as mass spectral similarity  
47 network or low library match score (level 2b), *in silico* based annotation (level 3), molecular formula  
48 match (level 4) and unknown spectral signals (level 5)<sup>3</sup>. A comprehensive dereplication may maximize  
49 annotation level 1 but involve a LC-MS/MS spectral library setup in identical analytical condition of  
50 matrix chemical profiling and is further limited to pure standards availability. Actually, authentic  
51 standard-centric annotation may identify only 1% to 10 % of all detected signals in a biological matrix  
52 but can be enriched using open mass spectral library (OMSL) resources to fill gaps with annotation  
53 level 2<sup>4</sup>.

54 Many OMSLs are freely available, such as GNPS, MassBank, MoNA, RIKEN, and HMDB<sup>5-8</sup> and  
55 immensely valuable for dereplication purposes. However, dealing with these resources is challenging  
56 due to the lack of standardized file formats and architecture. These libraries encompass a variety of file  
57 structures for mass spectral data, including ASCII-based formats like Mascot Generic Format (.MGF)  
58 and NIST MSP (.MSP), as well as MassBank records, JavaScript Object Notation (.JSON), Extensible  
59 Markup Language (.XML) or in the form of an SQLITE database<sup>9</sup>. While these formats generally follow  
60 a similar organizational schema—detailing compound spectra with core metadata on chemical  
61 identifiers (SMILES, INCHI, name, or adduct forms), experimental conditions (collision energy,  
62 ionization mode, polarity, or instrument type), and extended metadata for experimental measurements  
63 ( $m/z$  values, MS/MS fragments, and their intensities)—there is no uniformity in metadata field names,  
64 sequencing, or minimal requirements. This lack of standardization restricts OMSL compatibility with  
65 open-source processing software, making them prone to parsing and reading errors. For instance,  
66 OpenMS<sup>10</sup> only supports .MGF format while MS-DIAL<sup>11</sup> manages generic .MSP or MassBank records  
67 and MZMine<sup>12</sup> imports as .JSON, .MGF and .MSP files but may face parsing issues. Additionally, each  
68 OMSL favors a unique file format with its own metadata structure, based on undocumented and  
69 unversioned data models, limiting interoperability among LC-MS processing software and hindering  
70 the integrated use of multiple databases. MassBank is one of the few resources to offer guidelines  
71 describing these records based on a versioned repository (V2.6.0).

72 Recently, the Python package MatchMS<sup>13</sup> has proposed a pipeline to harmonize metadata and clean  
73 experimental values but focus mainly on data exploration using various MS/MS similarities measures.

74 For metadata enrichment related to chemical identifiers, another Python package MSMetaEnhancer  
75 have been added to MatchMS satellites tools<sup>14</sup>. Another shortcoming arises when using an OMSL:  
76 extracting a subset of interesting data proves difficult, given that most downloadable files are a  
77 concatenation of the two ionization modes, several collision energy methods, several instrument types,  
78 and a mix of predicted and experimental data. As a result, despite the great value of using one or several  
79 OMSLs, this appears challenging for dereplication of tandem mass spectra in daily work.

80 To bridge this gap, we introduce FragHub, a workflow that integrates diverse mass spectral libraries  
81 to streamline and enhance the annotation process. FragHUB support multiple OMSL formats (.MSP,  
82 .MGF, .JSON, .CSV, .XML) and harmonizes metadata using RDKit<sup>15</sup> and internal dictionaries. It allows  
83 for user-defined filtering options and handle outputs from MZMine's spectral library generation module,  
84 ensuring seamless integration of in-house databases. FragHub not only concatenates libraries from  
85 diverse sources into a unified format but also classifies the spectra according to chromatographic  
86 methods (GC/LC-MS), ionization modes (positive/negative), and data origin (predicted/experimental).  
87 Available as a Python package with a straightforward user interface, FragHub supports flexible  
88 parameter settings.

89 The processed libraries are compatible with Metabolomics data processing software such as MS-  
90 DIAL, MZMine3 or Flash Entropy Search<sup>16</sup>, but also interoperable with spectral data management  
91 software such as PeakForest<sup>17</sup>. A PeakForest instance for FragHub is accessible online, providing tools  
92 for viewing, browsing, and filtering spectral data through a web portal or API (available at  
93 <https://fraghub.peakforest.org/>).

94

## 95 **Materials and Methods**

96 FragHub's workflow was meticulously designed to parse and standardize spectral data across various  
97 formats, including .MSP, .MGF, .JSON, .CSV, and .XML, as derived from several widely utilized open  
98 mass spectral libraries. These operations involve detailed metadata normalization steps using RDKit,  
99 ensuring that data entries from disparate sources become interoperable. To validate and benchmark our  
100 approach, we utilized datasets encompassing over 790,000 spectra, demonstrating FragHub's ability to  
101 efficiently process and refine these entries for better usability in metabolic studies.

### 102 1. Open Mass Spectral Library Resources

103 The workflow was tested with a diversity of public software libraries (different data formats, diversity  
104 of metadata), four OMSLs were selected and downloaded in early January 2024 (see table 1).  
105 Additionally, an in-house database was created using MZMine3 to test outputs compatibility with  
106 FragHub. A step-by-step tutorial to create an in-house library is available in supplementary data. The  
107 dataset gathered for this work comprises 794,985 MS/MS spectra with the associated metadata.

108

### 109 **Tableau 1:OMSL list used to develop FragHub**

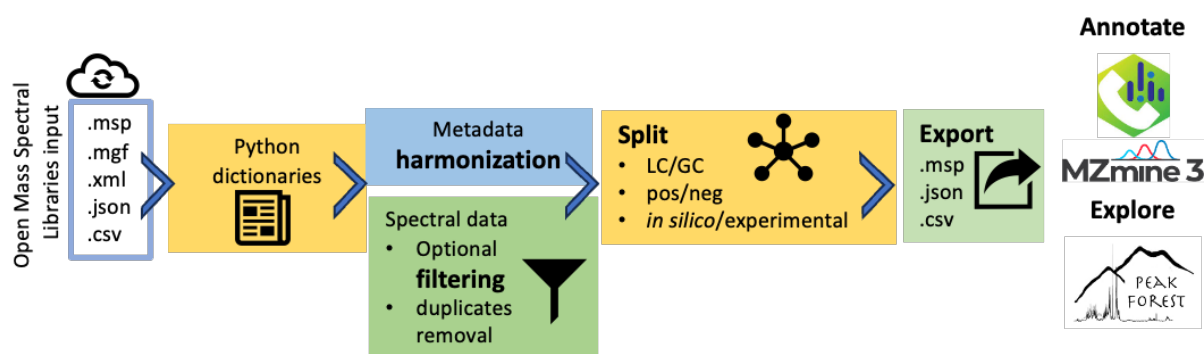
Spectral library name	URL	File format	Version	License	Spectra
MoNA	<a href="https://mona.fiehnlab.ucdavis.edu/downloads">https://mona.fiehnlab.ucdavis.edu/downloads</a>	.JSON	2024.01	CC-BY 4.0	190,359
MS-DIAL-VS17	<a href="http://prime.psc.riken.jp/comps/ms/msdial/main.html#MSP">http://prime.psc.riken.jp/comps/ms/msdial/main.html#MSP</a>	.MSP	2022.08	CC-BY 4.0	376,430
GNPS	<a href="https://gnps-external.ucsd.edu/gnpslibrary">https://gnps-external.ucsd.edu/gnpslibrary</a>	.MGF	2024.01 GNPS only	CC-0 1.0	63,935
MassBank	<a href="https://github.com/MassBank/MassBank-data">https://github.com/MassBank/MassBank-data</a>	.MSP	2023.11	CC-BY 4.0	164,261

110

111

112 2. FragHub Workflow

113



114

115 **Figure 1: FragHub workflow showing the 4 steps from OMSLs input to export files**

116

117 The initial step of the FragHub workflow involves parsing various data file formats, such as .MSP,  
 118 .MGF, .JSON, .CSV, and .XML, into field names and their corresponding values as delineated in Table  
 119 1. The workflow employs a mapping dictionary to translate current keys into standardized keys that  
 120 adhere to GNPS naming conventions, thereby ensuring compatibility with data reprocessing software  
 121 like MS-Dial, MZmine, and Flash Entropy Search which utilizes MSP and JSON for annotation.

122 To effectively manage duplicates and facilitate further data processing, FragHub generates a unique  
 123 hashing key (SHA-256) for each spectrum using the InChIKey and fragmentation spectra; if an  
 124 InChIKey is unavailable, the hashing key is derived from all available spectral data. This unique  
 125 identifier, termed 'FragHubID', simplifies the tracking and elimination of duplicate spectra both within  
 126 and across OMSLs. FragHubIDs are recorded in the “update.json” file, which helps in maintaining a

127 repository of processed spectra, ensuring that only new spectra are processed upon the addition of new  
128 OMSL entries, as configured by the user.

129 The workflow conducts a thorough cleaning and normalization of compound metadata and spectral  
130 data. It verifies the accuracy of SMILES, InChI, and InChIKey assignments, reallocating them as  
131 needed, and eliminates any spectra lacking both InChI and SMILES. RDKit is utilized to standardize  
132 chemical identifiers and calculate both exact and average molecular masses. Unparsable identifiers are  
133 removed, and any missing 'name' data are substituted with the corresponding molecule's InChI, where  
134 applicable. Non-specific values such as 'RT: 0.0' or 'adduct: unknown' are replaced with the placeholder  
135 "UNKNOWN". The workflow also updates adduct values, ion mode keys, and MS levels using a  
136 comprehensive mapping dictionary from the data directory, and tentatively calculates empty m/z  
137 precursor values based on the exact mass and identified adduct.

138 Instrument details (e.g., model types like QTOF or FT) and ionization modes (such as ESI or APCI)  
139 are normalized using the HUPO PSI mass spectrometry controlled vocabulary via an in-house  
140 hierarchical decision tree available in the data directory.

141 Spectra lacking essential information like SMILES, InChI, or a valid precursor m/z value, as well as  
142 those failing to meet user-specified filter criteria, are excluded. A detailed list of discarded spectra is  
143 compiled, highlighting the reasons for their removal.

144 Furthermore, FragHub annotates the 'predicted' field to distinguish between experimental and  
145 predicted spectra and normalizes retention times to minutes. Following metadata normalization, user-  
146 defined filters are applied through the graphical user interface to refine the peak list (Table S2).

147 Finally, the workflow segregates the spectra by ion detection mode (positive/negative), separation  
148 techniques (LC or GC), and categorizes them as experimental or predicted, removing any potential  
149 duplicates based on similar InChIKeys and their fragment lists. The entire process is efficiently  
150 completed in less than twenty minutes on a desktop computer equipped with an Intel Core i9-13900 and  
151 128 GB RAM DDR5, handling over a million spectra in various test formats.

152

### 153 3. OMSL benchmarking for annotation

154

155 In order to benchmark each OMSL for annotation purposes on a real dataset, raw data from Nicolle  
156 et al.<sup>18</sup> were used (<https://doi.org/10.5281/zenodo.8421008>). Quality control (pool of whole *Arabidopsis*  
157 *thaliana* extracts) and blank thermo .RAW data were imported into MS-Dial v5.231120.  
158 Chromatograms were deconvoluted, aligned using the same parameters as Nicolle et al. Then, filtered  
159 with the help of integrated MS-CleanR<sup>19</sup> with a blank ratio of 0.8; incorrect mass and ghost peak  
160 removed; a relative standard deviation of 40 and a relative mass defect between 50 and 3500. The  
161 alignment result was submitted to MS/MS based annotation using each OMSL processed by FragHub  
162 applying all default filters and exported in .MSP format. The following parameters were used for

163 spectral matches: Dot product score > 600; weighted dot product > 600; reverse dot product > 800;  
164 matched spectrum percentage > 25% and minimum number of matched peaks = 3.

165

166

167

#### 168 4. Chemical space representation

169

170 Chemical classes were deciphered using NPclassifier API<sup>20</sup>. PathwayNP and superclassNP were kept  
171 for each compound for figures coloration. The t-distributed Stochastic Neighbor Embedding (t-SNE)  
172 dimensionality reduction was calculated from PubChem fingerprints using a perplexity of 30 and an  
173 exaggeration of 1.

174

#### 175 5. PeakForest database

176

177 PeakForest is a multi-platform digital infrastructure for interoperable metabolite spectral data and  
178 metadata management. It captures and stores different types of metabolomics data from mass  
179 spectrometry and Nuclear magnetic resonance (NMR), providing users with valuable insights into  
180 metabolite identification and annotation processes. The infrastructure consists of a structured database,  
181 Application Programming Interfaces (API), a web interface and web services offering tools for  
182 browsing, managing and curating spectral data and metadata. Standardised procedures and formats have  
183 been implemented to guarantee information quality and interoperability. These features provide users  
184 with intuitive access to spectral data, facilitating efficient data annotation and analysis workflows.  
185 Finally, PeakForest is designed to facilitate the centralisation of data at laboratory level and to facilitate  
186 sharing between laboratories and public databases.

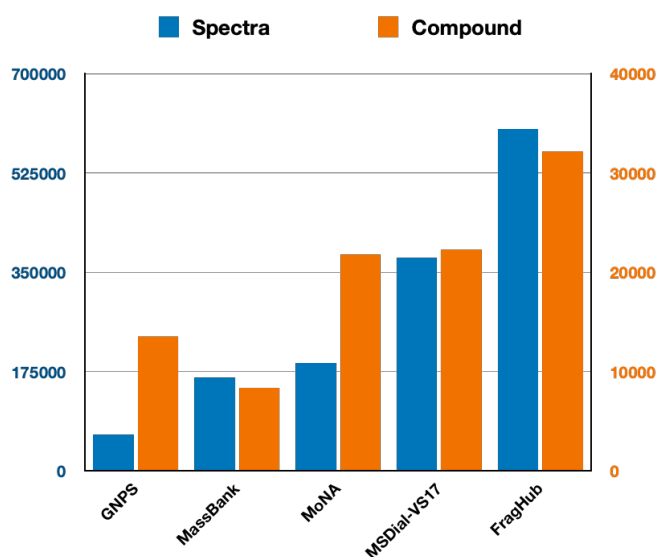
187

188

## 189 Results

190 FragHub, developed in the Python programming language, leverages four widely used open mass  
191 spectral libraries (OMSLs) for LC-MS-based metabolomic analysis. In this study, we specifically  
192 utilized GNPS-tagged databases in the .MGF format, comprising 13,507 compounds and 63,935  
193 spectra. Mona (MassBank of North America) significantly enriches our dataset with 21,839 unique  
194 compounds across 190,359 spectra, available in .MSP, .SDF, and .JSON formats. MassBank stands out  
195 for its spectral diversity, offering over 164,261 spectral datasets associated with 8,358 compounds.  
196 MSDial-VS17 represents a unique integration, merging several databases and in-house acquired spectra  
197 accounting for 376,430 spectra and 22,282 compounds. This dataset is the only library pre-split into  
198 positive ionization (PI) and negative ionization (NI) modes. For these latter two databases, the .MSP  
199 format has been utilized within FragHub. To showcase FragHub's adaptability, multiple formats were  
200 processed (as detailed in Table 1). The integration of these four OMSLs yields a combined total of  
201 794,985 spectra for 35,673 unique chemical identifiers. The FragHub data integration workflow refines  
202 this further to 602,744 spectra for 32,193 unique chemicals, as illustrated in Figure 2. Detailed logs of  
203 the spectra excluded during the OMSLs processing are maintained in Table S4.

204



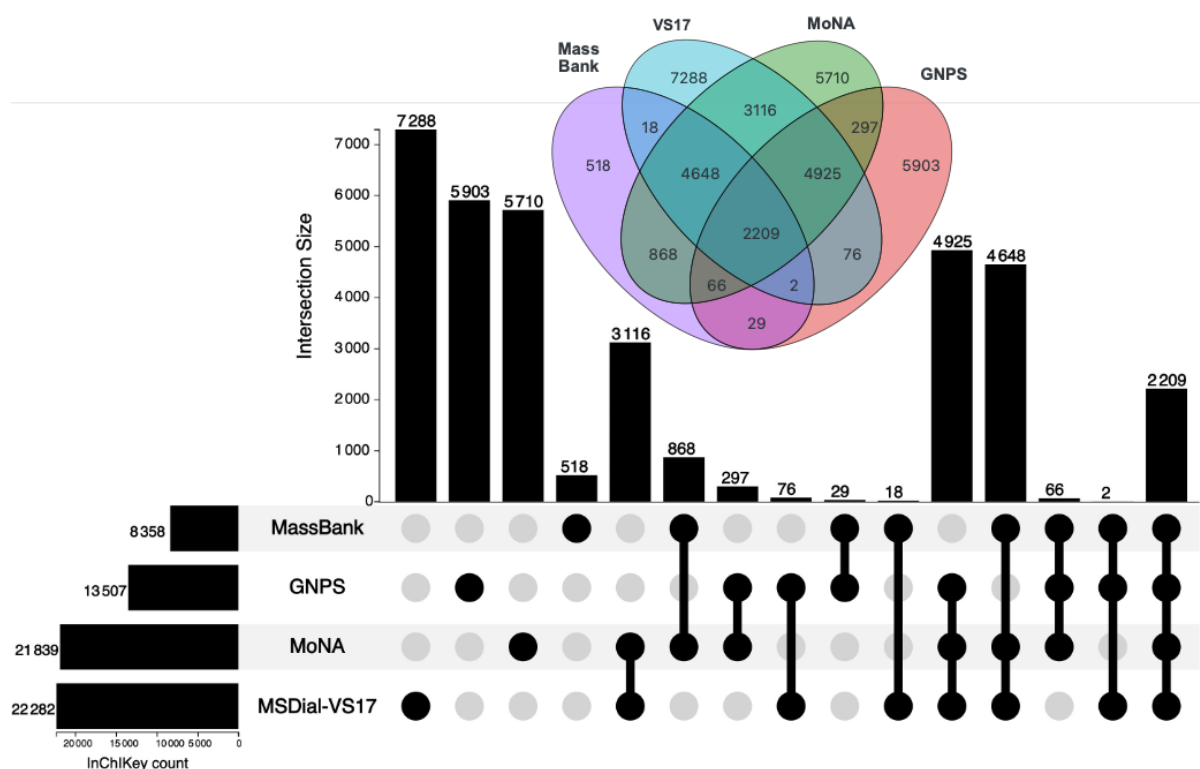
205

### 206 Figure 2: Integration Output Analysis

207 Bar plot displaying the counts of MS/MS spectra and unique InChIKeys derived from each OMSL. The left y-axis represents  
208 the number of spectra while the right y-axis shows the number of unique chemical identifiers. This visualization underscores  
209 the harmonization capabilities of FragHub, demonstrating its efficacy in integrating and deduplicating spectral data from  
210 diverse libraries.

211 Approximately 45% of chemicals are shared between two or more open mass spectral libraries  
212 (OMSLs), highlighting the interconnected nature of these resources. Conversely, 19,419 compounds  
213 are exclusive to a single OMSL. The FragHub workflow effectively reduces redundancy by eliminating  
214 about 200,000 duplicate spectra from an initial pool of 794,985, underscoring the diverse chemical  
215 compositions and experimental conditions—such as collision energy, instrument type, and adduct forms

216 of isolated pseudo-molecular ions—that characterize each library. The median number of spectra per  
217 compound ranges from 2 in GNPS to 12 in MassBank, illustrating significant spectral redundancy that  
218 can be tailored based on user preferences.  
219



220

221 **Figure 3: Compound Overlap among OMSLs**

222 Venn diagram and upset plot illustrating the intersection of unique compounds across various OMSLs. Each bar indicates the  
223 number of unique compounds exclusive to a single library or shared between multiple libraries, highlighting the  
224 complementary nature of the integrated libraries in covering broader chemical space.  
225

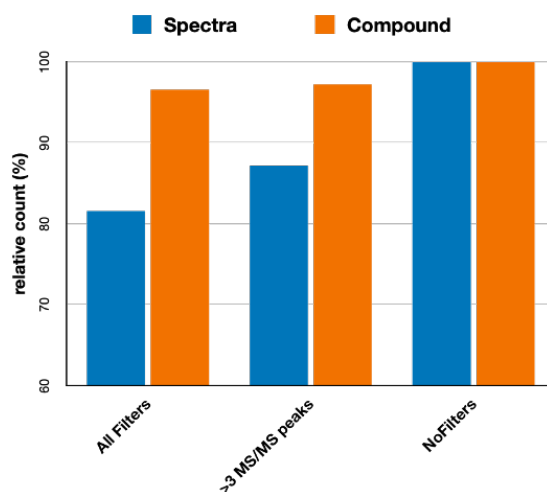
226

227 For example, applying a filter to remove spectra with fewer than three MS/MS signals results in a  
228 15% reduction in entries, as depicted in Figure 4. Further refinement is achieved through a second filter,  
229 which excludes spectra unless they meet a minimum threshold of three signals and two MS/MS peaks  
230 with intensities above 5%. This stringent criterion retains 81% of the total spectra while incurring a  
231 substantial loss of compounds, amounting to 3,5%, thereby optimizing the dataset for higher-quality  
232 annotations.  
233

232

233





234

235

236 **Figure 4: Filter Impact Analysis**

237 Bar plot quantifying the impact of applying default FragHub filters on the spectral and compound data retained from integrated

238 OMSLs. The plot compares the percentages of spectra and compounds retained with and without filtering, showcasing the

239 effectiveness of filters in enhancing data quality without significant loss of chemical diversity.

240

241 To assess the enhanced utility of integrated OMSLs for annotation tasks, we analyzed chemical

242 fingerprints from *Arabidopsis thaliana* using MS-Dial. The annotations were performed independently

243 on each OMSL as well as on the integrated dataset processed through the FragHub workflow. After

244 applying MS-CleanR filtration, a total of 435 features were detected in positive ionization mode. The

245 annotation process did not consider the retention time values and relied solely on accurate mass and

246 MS/MS fragmentation patterns.

247 The outcomes, depicted in Figure 5, demonstrate a direct relationship between the richness of the

248 compound library in each OMSL and the number of matches achieved: MassBank, with its 41 matches,

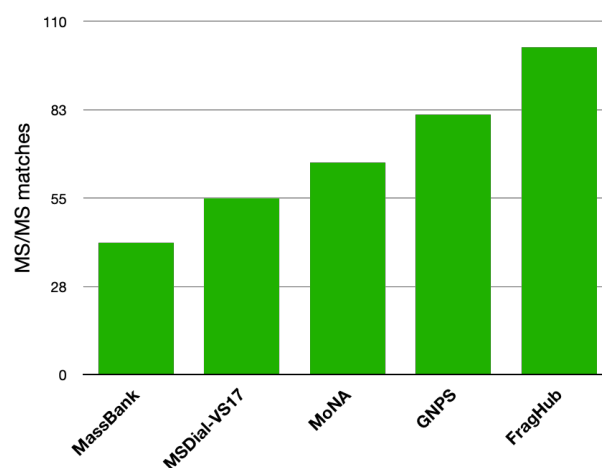
249 contrasts with the three other OMSLs, which, containing over 5,000 unique compounds each, yielded

250 between 55 and 81 matches. Remarkably, the consolidated file from FragHub, utilizing default filtering

251 criteria, successfully annotated 102 features, corresponding to 24% of the total detected features.

252

253



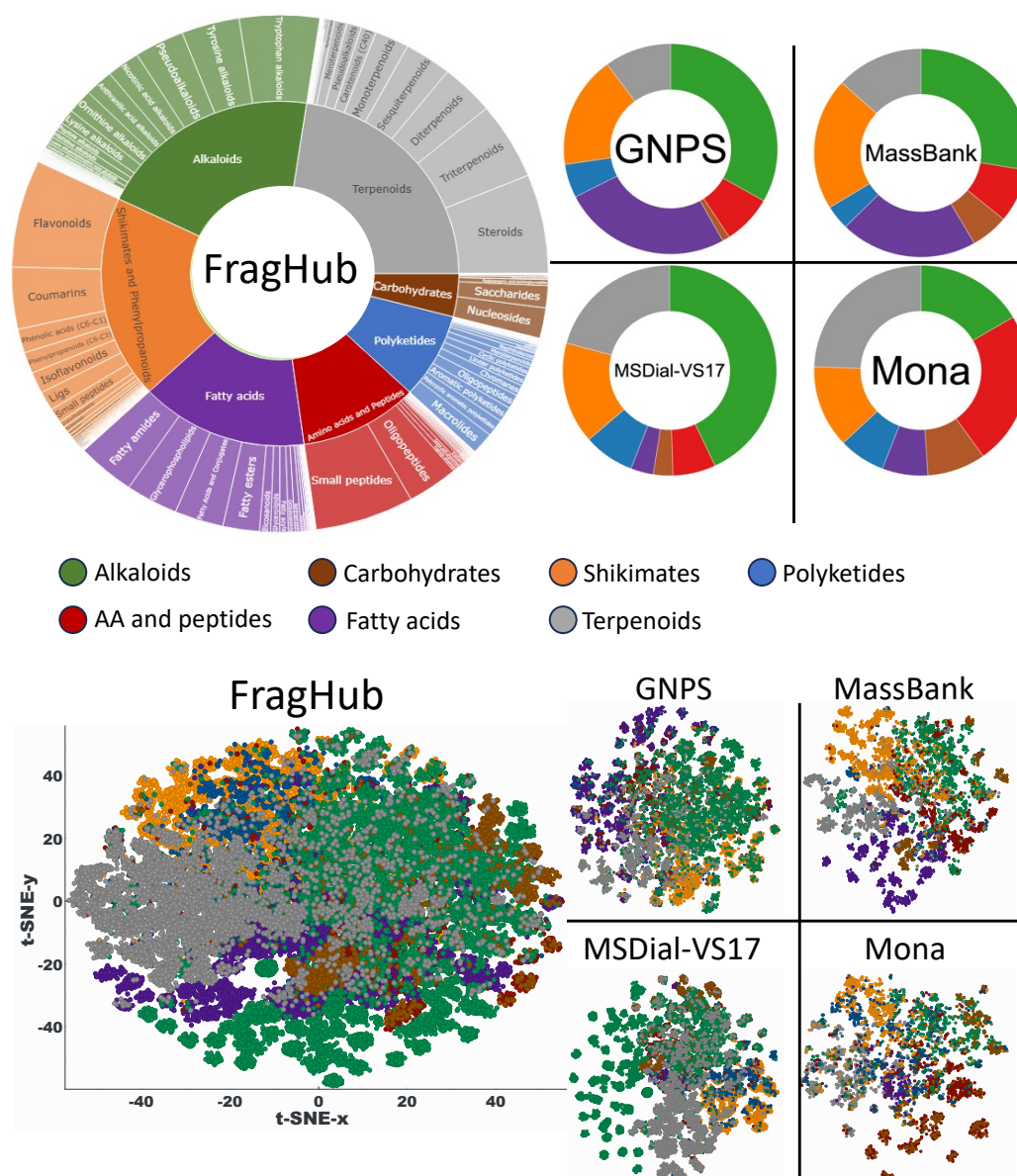
254

255 **Figure 5: Annotation Efficiency Comparison**

256 Bar plot showing the number of features successfully annotated from the Nicolle et al. dataset using individual and integrated  
 257 OMSLs under standard query conditions. This plot demonstrates the increased annotation capabilities achieved through the  
 258 integrated dataset, reflecting FragHub's enhancement for real dataset annotation.

259

260 The distribution of chemical classes across each OMSL highlights the unique chemical diversity they  
 261 cover. Fatty acids predominate in GNPS and MassBank, whereas alkaloids are prominently featured in  
 262 GNPS and MSDial-VS17. Mona is rich in carbohydrates, amino acids, and peptides. In contrast,  
 263 shikimates, phenylpropanoids, and terpenoids are more evenly distributed across the OMSLs, as shown  
 264 in the top panel of Figure 6. The chemical space of each OMSL was analyzed using a t-SNE  
 265 dimensionality reduction approach based on PubChem fingerprints. This method effectively reveals  
 266 local clusters and the overall spatial distribution of compounds, facilitating an intuitive visualization of  
 267 how different chemical classes aggregate. Typically, compounds within the same class cluster together,  
 268 with each class occupying distinct regions in the t-SNE plot. GNPS and MSDial show denser  
 269 distributions, particularly in the areas representing terpenoids and alkaloids, whereas Mona spans a  
 270 broader area for carbohydrates, and MassBank is extensively spread across regions rich in shikimates  
 271 and phenylpropanoids. Collectively, the integration of these OMSLs through FragHub achieves a  
 272 comprehensive and dense coverage of chemical space across all compound classes.



273

274 **Figure 6: Chemical Space Coverage**

275 t-SNE plots overlaid with donut charts depicting the distribution of metabolite classes within each OMSL and the integrated  
 276 dataset. The t-SNE plots provide a two-dimensional representation of the chemical spaces covered by each library, with colors  
 277 indicating different chemical classes based on NPClassifier ontology. The donut charts further detail the proportion of each  
 278 metabolite class, illustrating the enriched diversity achieved through data integration.

279 **Discussion**

280

281 The growing number of publications in metabolomics underscores its significance within the omics  
 282 landscape, yet the relevance of the results stemming from this approach is largely dependent on the  
 283 quality of annotations derived from spectrometric signals<sup>21</sup>. In this context, OMSLs are key for  
 284 supporting experimental spectral matching and enhancing annotation rate from untargeted LC-MS  
 285 fingerprints. The aim of FragHub workflow is to optimize the use of OMSLs for end-users in the field  
 286 of untargeted LC-MS based metabolomic. Four OMSLs have been used in various formats to  
 287 demonstrate the FragHub integration pipeline (Figure 2). A primary challenge in this integration was

288 the normalization of data fields and values from diverse sources. For example, we identified ten distinct  
289 keys for ionization states and normalized 487 instrument names and 154 adduct descriptions to 307 and  
290 111, respectively, as detailed in Table S1. The harmonization of collision energies was not addressed  
291 due to their varied and non-standardized measures (around 70 different formats), highlighting the  
292 critical need for standardized data practices as recommended by MassBank.as recommended in the  
293 MassBank documentation for instance<sup>22</sup>.

294 Approximately 50% of unique compounds and 20% of spectral duplicates were observed across the  
295 OMSLs, indicating that while high redundancy can improve annotation rates, it might also lead to  
296 inconsistencies, particularly when using dot product and reverse dot product scoring systems that are  
297 highly sensitive to fragment number and intensity. To mitigate these issues, FragHub implements filters  
298 that maintain data integrity without compromising compound diversity, as shown in Figure 4.  
299 Furthermore, MS/MS data denoising may be applied by plugging FragHub outputs to Libgen<sup>23</sup> or  
300 alternative scoring approach<sup>24</sup>.

301 The integration of OMSLs used here significantly expands the compound diversity and chemical  
302 space coverage and increase annotation rate of untargeted chemical profiling (Figure 5 and 6). In the  
303 context of holistic approaches, deciphering the interplay of metabolome dynamics across organisms or  
304 environments is challenging. The use of large mass spectral libraries extends metabolome coverage  
305 outside of expected results enabling a comprehensive understanding of complex systems. We measured  
306 32,193 unique compounds after OMSLs integration which is rather low compared to the diversity of  
307 natural products estimated to be several million molecules<sup>25</sup>. Moreover, the chemical classes covered  
308 by OMSLs contrast with the distribution of natural products databases such as the Dictionary of Natural  
309 Products with an over-representation of alkaloids and polypeptides in OMSL, while terpenoids and fatty  
310 acids represent the most diverse group in natural product catalogues<sup>26</sup>. This disparity underscores the  
311 necessity for orthogonal strategies to fill this gap like raw data digging of mass spectral similarity  
312 networks<sup>27</sup> or *in silico* MS/MS prediction tools based on chemical identifiers<sup>28</sup>. The FragHub integration  
313 workflow may help to organize data and explore fragmentation mechanism behavior to set up training  
314 sets for deep learning-based strategies.

315 The FragHub code can handle various input formats and has been multithreaded to process  
316 approximately 100,000 spectra per minute (table S3) which allows the integration of large OMSLs in  
317 reasonable time on a personal computer. A simple graphical user interface enables users to select  
318 filtering options and data format outputs using distinct profiles. This allows shaping scenarios for  
319 specific needs such as in-house database handling or simple .CSV outputs to analyze OMSLs, then filter  
320 on specific metadata (e.g., instrument type) and reintegration in .MSP or .JSON formats for instance.

321 To demonstrate the potential of this data standardization and structuring work, the compounds and  
322 their LC-MSMS spectra were also imported and stored in a dedicated PeakForest database. The web  
323 application provided enables users, for example, to browse and search for specific chemical names or

324 spectral metadata. It also provides a REST web service to support massive queries submitted by third-  
325 party software or bioinformatics pipelines for metabolomics data annotation. PeakForest has been  
326 initially developed to store and manage high-quality spectral data in terms of metadata. The FragHub  
327 instance of PeakForest can be used to put online a collection of sub-banks in MSP format, compiled for  
328 example by instrument type. By exploiting the various resources made available by the community and  
329 used in the FragHub pipeline, we were able to compile a very large number of MSMS spectra. This  
330 work once again highlights the need to open up more and more new spectral data, acquired on recent  
331 instruments and supplemented with rich, controlled metadata, in order to increase annotation coverage  
332 of LC-MS fingerprints.

333 The integration of multiple mass spectral libraries through FragHub represents a significant advance  
334 in the metabolomics field, facilitating a deeper understanding of metabolite environments through  
335 enhanced data quality and accessibility. Moreover, FragHub's flexible architecture allows for the rapid  
336 incorporation of new data sources, which is critical given the rapid evolution of mass spectrometry  
337 libraries. By addressing the critical challenges of data standardization and compatibility, FragHub  
338 provides researchers with powerful tools to unlock the full potential of metabolomic studies.

339

340

## 341 **References**

342

343 (1) Wolfender, J.-L.; Marti, G.; Thomas, A.; Bertrand, S. Current Approaches and Challenges for the Metabolite Profiling of  
344 Complex Natural Extracts. *J. Chromatogr. A* 2015, 1382, 136–164. <https://doi.org/10.1016/j.chroma.2014.10.091>.

345 (2) Sumner, L. W.; Amberg, A.; Barrett, D.; Beale, M. H.; Beger, R.; Daykin, C. A.; Fan, T. W.-M.; Fiehn, O.; Goodacre, R.;  
346 Griffin, J. L.; Hankemeier, T.; Hardy, N.; Harnly, J.; Higashi, R.; Kopka, J.; Lane, A. N.; Lindon, J. C.; Marriott, P.; Nicholls, A. W.;  
347 Reily, M. D.; Thaden, J. J.; Viant, M. R. Proposed Minimum Reporting Standards for Chemical Analysis: Chemical Analysis Working  
348 Group (CAWG) Metabolomics Standards Initiative (MSI). *Metabolomics* 2007, 3 (3), 211–221. [https://doi.org/10.1007/s11306-007-](https://doi.org/10.1007/s11306-007-0082-2)  
349 0082-2.

350 (3) Schymanski, E. L.; Jeon, J.; Gulde, R.; Fenner, K.; Ruff, M.; Singer, H. P.; Hollender, J. Identifying Small Molecules via  
351 High Resolution Mass Spectrometry: Communicating Confidence. *Environ. Sci. Technol.* 2014, 48 (4), 2097–2098.  
352 <https://doi.org/10.1021/es5002105>.

353 (4) Tsugawa, H.; Rai, A.; Saito, K.; Nakabayashi, R. Metabolomics and Complementary Techniques to Investigate the Plant  
354 Phytochemical Cosmos. *Nat. Prod. Rep.* 2021, 10.1039/D1NP00014D. <https://doi.org/10.1039/D1NP00014D>.

355 (5) Aron, A. T.; Gentry, E.; McPhail, K. L.; Nothias, L. F.; Nothias-Esposito, M.; Boulimani, A.; Petras, D.; Gauglitz, J. M.;  
356 Sikora, N.; Vargas, F.; J. van der Hooft, J. J.; Ernst, M.; Kang, K. B.; Aceves, C. M.; Caraballo-Rodríguez, A. M.; Koester, I.; Weldon,  
357 K. C.; Bertrand, S.; Roullier, C.; Sun, K.; Tehan, R. M.; Boya, C. A.; Martin, H., C.; Gutiérrez, M.; Ulloa, A. M.; Mora, J. A. T.;  
358 Mojica-Flores, R.; Lakey-Betitia, J.; Vázquez-Chaves, V.; I. Calderón, A.; Tayler, N.; Keyzers, R. A.; Tugizimana, F.; Ndlovu, N.;  
359 Aksenov, A. A.; Jarmusch, A. K.; Schmid, R.; Truman, A. W.; Bandeira, N.; Wang, M.; Dorrestein, P. Reproducible Molecular  
360 Networking Of Untargeted Mass Spectrometry Data Using GNPS.; preprint; 2019. <https://doi.org/10.26434/chemrxiv.9333212.v1>.

361 (6) Horai, H.; Arita, M.; Kanaya, S.; Nihei, Y.; Ikeda, T.; Suwa, K.; Ojima, Y.; Tanaka, K.; Tanaka, S.; Aoshima, K.; Oda, Y.;  
362 Kakazu, Y.; Kusano, M.; Tohge, T.; Matsuda, F.; Sawada, Y.; Hirai, M. Y.; Nakanishi, H.; Ikeda, K.; Akimoto, N.; Maoka, T.;  
363 Takahashi, H.; Ara, T.; Sakurai, N.; Suzuki, H.; Shibata, D.; Neumann, S.; Iida, T.; Tanaka, K.; Funatsu, K.; Matsuura, F.; Soga, T.;  
364 Taguchi, R.; Saito, K.; Nishioka, T. MassBank: A Public Repository for Sharing Mass Spectral Data for Life Sciences. *J. Mass*  
365 *Spectrom.* 2010, 45 (7), 703–714. <https://doi.org/10.1002/jms.1777>.

- 366 (7) Tsugawa, H.; Ikeda, K.; Takahashi, M.; Satoh, A.; Mori, Y.; Uchino, H.; Okahashi, N.; Yamada, Y.; Tada, I.; Bonini, P.;  
367 Higashi, Y.; Okazaki, Y.; Zhou, Z.; Zhu, Z.-J.; Koelmel, J.; Cajka, T.; Fiehn, O.; Saito, K.; Arita, M.; Arita, M. A Lipidome Atlas in  
368 MS-DIAL 4. *Nat. Biotechnol.* 2020, 1–5. <https://doi.org/10.1038/s41587-020-0531-2>.
- 369 (8) Wishart, D. S.; Guo, A.; Oler, E.; Wang, F.; Anjum, A.; Peters, H.; Dizon, R.; Sayeeda, Z.; Tian, S.; Lee, B. L.; Berjanskii,  
370 M.; Mah, R.; Yamamoto, M.; Jovel, J.; Torres-Calzada, C.; Hiebert-Giesbrecht, M.; Lui, V. W.; Varshavi, D.; Varshavi, D.; Allen,  
371 D.; Arndt, D.; Khetarpal, N.; Sivakumaran, A.; Harford, K.; Sanford, S.; Yee, K.; Cao, X.; Budinski, Z.; Liigand, J.; Zhang, L.; Zheng,  
372 J.; Mandal, R.; Karu, N.; Dambrova, M.; Schiöth, H. B.; Greiner, R.; Gautam, V. HMDB 5.0: The Human Metabolome Database for  
373 2022. *Nucleic Acids Res.* 2022, 50 (D1), D622–D631. <https://doi.org/10.1093/nar/gkab1062>.
- 374 (9) Kind, T.; Tsugawa, H.; Cajka, T.; Ma, Y.; Lai, Z.; Mehta, S. S.; Wohlgemuth, G.; Barupal, D. K.; Showalter, M. R.; Arita,  
375 M.; Fiehn, O. Identification of Small Molecules Using Accurate Mass MS/MS Search. *Mass Spectrom. Rev.* 2018, 37 (4), 513–532.  
376 <https://doi.org/10.1002/mas.21535>.
- 377 (10) Röst, H. L.; Sachsenberg, T.; Aiche, S.; Bielow, C.; Weisser, H.; Aicheler, F.; Andreotti, S.; Ehrlich, H.-C.; Gutenbrunner,  
378 P.; Kenar, E.; Liang, X.; Nahnsen, S.; Nilse, L.; Pfeuffer, J.; Rosenberger, G.; Rurik, M.; Schmitt, U.; Veit, J.; Walzer, M.; Wojnar,  
379 D.; Wolski, W. E.; Schilling, O.; Choudhary, J. S.; Malmström, L.; Aebersold, R.; Reinert, K.; Kohlbacher, O. OpenMS: A Flexible  
380 Open-Source Software Platform for Mass Spectrometry Data Analysis. *Nat. Methods* 2016, 13 (9), 741–748.  
381 <https://doi.org/10.1038/nmeth.3959>.
- 382 (11) Tsugawa, H.; Cajka, T.; Kind, T.; Ma, Y.; Higgins, B.; Ikeda, K.; Kanazawa, M.; VanderGheynst, J.; Fiehn, O.; Arita, M.  
383 MS-DIAL: Data-Independent MS/MS Deconvolution for Comprehensive Metabolome Analysis. *Nat. Methods* 2015, 12 (6), 523–  
384 526. <https://doi.org/10.1038/nmeth.3393>.
- 385 (12) Schmid, R.; Heuckeroth, S.; Korf, A.; Smirnov, A.; Myers, O.; Dyrlund, T. S.; Bushuiev, R.; Murray, K. J.; Hoffmann, N.;  
386 Lu, M.; Sarvepalli, A.; Zhang, Z.; Fleischauer, M.; Dührkop, K.; Wesner, M.; Hoogstra, S. J.; Rudt, E.; Mokshyna, O.; Brungs, C.;  
387 Ponomarov, K.; Mutabdzija, L.; Damiani, T.; Pudney, C. J.; Earll, M.; Helmer, P. O.; Fallon, T. R.; Schulze, T.; Rivas-Ubach, A.;  
388 Bilbao, A.; Richter, H.; Nothias, L.-F.; Wang, M.; Orešič, M.; Weng, J.-K.; Böcker, S.; Jeibmann, A.; Hayen, H.; Karst, U.; Dorrestein,  
389 P. C.; Petras, D.; Du, X.; Pluskal, T. Integrative Analysis of Multimodal Mass Spectrometry Data in MZmine 3. *Nat. Biotechnol.*  
390 2023, 41 (4), 447–449. <https://doi.org/10.1038/s41587-023-01690-2>.
- 391 (13) Huber, F.; Verhoeven, S.; Meijer, C.; Spreuw, H.; Castilla, E.; Geng, C.; van der Hooft, J.; Rogers, S.; Belloum, A.;  
392 Diblen, F.; Spaaks, J. Matchms - Processing and Similarity Evaluation of Mass Spectrometry Data. *J. Open Source Softw.* 2020, 5  
393 (52), 2411. <https://doi.org/10.21105/joss.02411>.
- 394 (14) Troják, M.; Hecht, H.; Čech, M.; Price, E. J. MSMetaEnhancer: A Python Package for Mass Spectrametadata Annotation.  
395 *J. Open Source Softw.* 2022, 7 (79), 4494. <https://doi.org/10.21105/joss.04494>.
- 396 (15) Landrum, G. Rdkit Documentation. Release 2013, 1 (1–79), 4.
- 397 (16) Li, Y.; Fiehn, O. Flash Entropy Search to Query All Mass Spectral Libraries in Real Time. *Nat. Methods* 2023, 20 (10),  
398 1475–1478. <https://doi.org/10.1038/s41592-023-02012-9>.
- 399 (17) Paulhe, N.; Canlet, C.; Damont, A.; Peyriga, L.; Durand, S.; Deborde, C.; Alves, S.; Bernillon, S.; Berton, T.; Bir, R.;  
400 Bouville, A.; Cahoreau, E.; Centeno, D.; Costantino, R.; Debrauwer, L.; Delabrière, A.; Duperier, C.; Emery, S.; Flandin, A.;  
401 Hohenester, U.; Jacob, D.; Joly, C.; Jousse, C.; Lagree, M.; Lamari, N.; Lefebvre, M.; Lopez-Piffet, C.; Lyan, B.; Maucourt, M.;  
402 Migne, C.; Olivier, M.-F.; Rathahao-Paris, E.; Petriacq, P.; Pinelli, J.; Roch, L.; Roger, P.; Roques, S.; Tabet, J.-C.; Tremblay-Franco,  
403 M.; Traïkia, M.; Warnet, A.; Zhendre, V.; Rolin, D.; Jourdan, F.; Thévenot, E.; Moing, A.; Jamin, E.; Fenaille, F.; Junot, C.; Pujos-  
404 Guillot, E.; Giacomoni, F. PeakForest: A Multi-Platform Digital Infrastructure for Interoperable Metabolite Spectral Data and  
405 Metadata Management. *Metabolomics* 2022, 18 (6), 40. <https://doi.org/10.1007/s11306-022-01899-3>.
- 406 (18) Nicolle, C.; Gayrard, D.; Noël, A.; Hortala, M.; Amiel, A.; Grat-Simeone, S.; Le Ru, A.; Marti, G.; Pernodet, J.-L.; Lautru,  
407 S.; Dumas, B.; Rey, T. Root Associated Streptomyces Produce Galbonolides to Modulate Plant Immunity and Promote Rhizosphere  
408 Colonisation. 2024. <https://doi.org/10.1101/2024.01.20.576418>.
- 409 (19) Fraisier-Vannier, O.; Chervin, J.; Cabanac, G.; Puech, V.; Fournier, S.; Durand, V.; Amiel, A.; André, O.; Benamar, O. A.;  
410 Dumas, B.; Tsugawa, H.; Marti, G. MS-CleanR: A Feature-Filtering Workflow for Untargeted LC–MS Based Metabolomics. *Anal.*  
411 *Chem.* 2020, 92 (14), 9971–9981. <https://doi.org/10.1021/acs.analchem.0c01594>.

412 (20) Kim, H. W.; Wang, M.; Leber, C. A.; Nothias, L.-F.; Reher, R.; Kang, K. B.; van der Hooft, J. J. J.; Dorrestein, P. C.;  
413 Gerwick, W. H.; Cottrell, G. W. NPClassifier: A Deep Neural Network-Based Structural Classification Tool for Natural Products. *J.*  
414 *Nat. Prod.* 2021, 84 (11), 2795–2807. <https://doi.org/10.1021/acs.jnatprod.1c00399>.

415 (21) Theodoridis, G.; Gika, H.; Raftery, D.; Goodacre, R.; Plumb, R. S.; Wilson, I. D. Ensuring Fact-Based Metabolite  
416 Identification in Liquid Chromatography–Mass Spectrometry-Based Metabolomics. *Anal. Chem.* 2023, 95 (8), 3909–3916.  
417 <https://doi.org/10.1021/acs.analchem.2c05192>.

418 (22) MassBank Documentation. [https://massbank.github.io/MassBank-documentation/contributor\\_documentation.html](https://massbank.github.io/MassBank-documentation/contributor_documentation.html).

419 (23) Kong, F.; Keshet, U.; Shen, T.; Rodriguez, E.; Fiehn, O. LibGen: Generating High Quality Spectral Libraries of Natural  
420 Products for EAD-, UVPD-, and HCD-High Resolution Mass Spectrometers. *Anal. Chem.* 2023, 95 (46), 16810–16818.  
421 <https://doi.org/10.1021/acs.analchem.3c02263>.

422 (24) Li, Y.; Kind, T.; Folz, J.; Vaniya, A.; Mehta, S. S.; Fiehn, O. Spectral Entropy Outperforms MS/MS Dot Product Similarity  
423 for Small-Molecule Compound Identification. *Nat. Methods* 2021, 18 (12), 1524–1531. <https://doi.org/10.1038/s41592-021-01331-z>.

424 (25) Medema, M. H.; de Rond, T.; Moore, B. S. Mining Genomes to Illuminate the Specialized Chemistry of Life. *Nat. Rev.*  
425 *Genet.* 2021, 22 (9), 553–571. <https://doi.org/10.1038/s41576-021-00363-7>.

426 (26) Chassagne, F.; Cabanac, G.; Hubert, G.; David, B.; Marti, G. The Landscape of Natural Product Diversity and Their  
427 Pharmacological Relevance from a Focus on the Dictionary of Natural Products®. *Phytochem. Rev.* 2019, 18 (3), 601–622.  
428 <https://doi.org/10.1007/s11101-019-09606-2>.

429 (27) Bittremieux, W.; Avalon, N. E.; Thomas, S. P.; Kakhkhorov, S. A.; Aksenov, A. A.; Gomes, P. W. P.; Aceves, C. M.;  
430 Caraballo-Rodríguez, A. M.; Gauglitz, J. M.; Gerwick, W. H.; Huan, T.; Jarmusch, A. K.; Kaddurah-Daouk, R. F.; Kang, K. B.; Kim,  
431 H. W.; Kondić, T.; Mannochio-Russo, H.; Meehan, M. J.; Melnik, A. V.; Nothias, L.-F.; O'Donovan, C.; Panitchpakdi, M.; Petras,  
432 D.; Schmid, R.; Schymanski, E. L.; Van Der Hooft, J. J. J.; Weldon, K. C.; Yang, H.; Xing, S.; Zemlin, J.; Wang, M.; Dorrestein, P.  
433 C. Open Access Repository-Scale Propagated Nearest Neighbor Suspect Spectral Library for Untargeted Metabolomics. *Nat.*  
434 *Commun.* 2023, 14 (1), 8488. <https://doi.org/10.1038/s41467-023-44035-y>.

435 (28) Mallowney, M. W.; Duncan, K. R.; Elsayed, S. S.; Garg, N.; van der Hooft, J. J. J.; Martin, N. I.; Meijer, D.; Terlouw, B.  
436 R.; Biermann, F.; Blin, K.; Durairaj, J.; Gorostiola González, M.; Helfrich, E. J. N.; Huber, F.; Leopold-Messer, S.; Rajan, K.; de  
437 Rond, T.; van Santen, J. A.; Sorokina, M.; Balunas, M. J.; Beniddir, M. A.; van Bergeijk, D. A.; Carroll, L. M.; Clark, C. M.; Clevert,  
438 D.-A.; Dejong, C. A.; Du, C.; Ferrinho, S.; Grisoni, F.; Hofstetter, A.; Jespers, W.; Kalinina, O. V.; Kautsar, S. A.; Kim, H.; Leao, T.  
439 F.; Masschelein, J.; Rees, E. R.; Reher, R.; Reker, D.; Schwaller, P.; Segler, M.; Skinnider, M. A.; Walker, A. S.; Willighagen, E. L.;  
440 Zdrzil, B.; Ziemert, N.; Goss, R. J. M.; Guyomard, P.; Volkamer, A.; Gerwick, W. H.; Kim, H. U.; Müller, R.; van Wezel, G. P.; van  
441 Westen, G. J. P.; Hirsch, A. K. H.; Linington, R. G.; Robinson, S. L.; Medema, M. H. Artificial Intelligence for Natural Product Drug  
442 Discovery. *Nat. Rev. Drug Discov.* 2023, 22 (11), 895–916. <https://doi.org/10.1038/s41573-023-00774-7>.

443

## 444 SUPPORTING INFORMATION

- 445 • Supplementary Tables comprising table S1 to S4 in .PDF
- 446 • Tutorial for FragHub installation and usage in .PDF
- 447 • Tutorial to set-up in-house library using MZMine in .PDF

448

## 449 AVAILABILITY

450 FragHub code can be forked, cloned or downloaded on GitHub at the following address:  
451 <https://github.com/eMetaboHUB/FragHub>.

452 FragHub is available with a pre-built data structure to facilitate the end-user processing. A tutorial is available on GitHub  
453 repository and in supplementary data.

454 OMSLs processed in this study are available on Zenodo repository: <https://doi.org/10.5281/zenodo.11057687>.

455

## 456 AUTHOR INFORMATION

### 457 Corresponding Author

15

458 \* Guillaume Marti, [guillaume.marti@univ-tlse3.fr](mailto:guillaume.marti@univ-tlse3.fr) Laboratoire de Recherche en Sciences Végétales, Metatoul-AgromiX  
459 Platform, Université de Toulouse, CNRS, INP, 24 Chemin de Borde Rouge, Auzeville, 31320, Auzeville-Tolosane, France.

460

#### 461 **Acknowledgment**

462 The FragHub and PeakForest project is supported by the French National Facility in Metabolomics & Fluxomics, MetaboHUB  
463 (11-INBS-0010), launched by the French Ministry of Research and Higher Education and the French ANR funding agency  
464 within the Programme “France 2030”. We also acknowledge all contributors to open mass spectral libraries.

465

#### 466 **Author Contributions**

467 GM proposed the study; AD and SH developed the python package; EJ and BL setup values dictionaries; NP and FG  
468 developed the peakforest instance; GM, GC and YG benchmarked and reviewed the workflow. The manuscript was written  
469 through contributions of all authors and all authors have given approval to the final version of the manuscript.

470