

1 **Generalized tree structure to annotate untargeted metabolomics and stable isotope** 2 **tracing data**

3

4 Shuzhao Li* and Shujian Zheng

5 Jackson Laboratory for Genomic Medicine, Farmington, CT 06032, USA

6 *Corresponding author, E-mail: shuzhao.li@jax.org

7

8 **Abstract**

9 In untargeted metabolomics, multiple ions are often measured for each original metabolite,
10 including isotopic forms and in-source modifications, such as adducts and fragments. Without
11 prior knowledge of the chemical identity or formula, computational organization and interpretation
12 of these ions is challenging, which is the deficit of previous software tools that perform the task
13 using network algorithms. We propose here a generalized tree structure to annotate ions to
14 relationships to the original compound and infer neutral mass. An algorithm is presented to
15 convert mass distance networks to this tree structure with high fidelity. This method is useful for
16 both regular untargeted metabolomics and stable isotope tracing experiments. It is implemented
17 as a Python package (khipu), and provides a JSON format for easy data exchange and software
18 interoperability. By generalized pre-annotation, khipu makes it feasible to connect metabolomics
19 data with common data science tools, and supports flexible experimental designs.

20

21 **Introduction**

22 Metabolomics is becoming an increasingly important tool to biomedicine. Untargeted LC-MS
23 (liquid chromatography-mass spectrometry) metabolomics is key to perform high-coverage
24 chemical analysis and discoveries. The term "annotation" in metabolomics often includes i) the
25 assignment of measured ions to their original compounds, and ii) establishing the identity of the
26 compounds (Domingo-Almenara et al, 2018; Blaženović et al, 2019). For clarity, we refer the first
27 step as "pre-annotation" in this paper, which is the assignment of isotopes, adducts and fragments
28 to the unique compounds. Correct pre-annotation will greatly facilitate the later step of
29 identification, by reducing errors on analyzing and searching the redundant ions. Multiple software
30 tools have been developed for this purpose of pre-annotation, including CAMERA (Kuhl et al,
31 2012), Mz.unity (Mahieu et al, 2016), xMSannotator (Uppal et al, 2017), MS-FLO (DeFelice et al,
32 2017), MetNet (Naake and Fernie, 2018), CliqueMS (Senan et al, 2019), Binner (Kachman et al,
33 2020) and NetID (Chen et al, 2021).

34 In high-resolution mass spectrometry, the m/z (mass to charge ratio) difference between isotopes
35 is usually resolved unambiguously. Adducts are formed in the ionization process, therefore, those
36 from the same original compound should have the same retention time in chromatography.
37 Besides adduct ions, formation of conjugates and fragments (including neutral loss) also belongs
38 to in-source modifications. Isotopes, adducts and fragments are often referred as redundant or
39 degenerate peaks in LC-MS literature. All pre-annotation tools utilize the m/z differences between
40 peaks, which correspond to the mass differential between isotopes, or between atoms or chemical
41 groups. In addition, having the same retention time is a critical requirement to group these
42 degenerate peaks. Some tools also use similarity in the shape of elution peaks and sometimes
43 statistical correlation between peak intensity across samples. Such correlations can be supporting
44 evidence but are not a prerequisite (Mahieu et al, 2016).

45
46 Most pre-annotation tools use a network representation of degenerate peaks. Because the
47 pairwise relationships between peaks are established first, then it is natural to connect the pairs
48 into networks by using pairwise relationships as edges and shared peaks as nodes. Such
49 networks still contain redundant and often erroneous edges. The main challenge remains to
50 resolve how all peaks are generated from the same original compound, which requires a) inferring
51 the neutral mass of the original compound, and b) establishing the relationship of all peaks to the
52 original compound. Given the difficulty of organizing this information in untargeted metabolomics,
53 the coverage of untargeted analyses is often called to question.

54
55 A couple of notable studies tried to address the question of coverage using isotope tracing in
56 untargeted metabolomics, and suggested that a small number of metabolites are actually
57 measured and the majority of peaks are "junks", either from contaminations, isotopes or LC-MS
58 artifacts (Mahieu and Patti, 2017; Wang et al, 2019). A new challenge also arose that analyzing
59 these isotope tracing data by global metabolomic is not trivial. So far, isotope tracing experiments
60 usually require targeted metabolites and specialized software (Chokkathukalam et al, 2013;
61 Bueschl et al, 2017; Previs and Downes, 2020; Rahim et al, 2022). In untargeted analysis, without
62 prior knowledge of the chemical formulas, special experimental designs are required and the
63 software tools are tied to the designs, which are the cases for $X^{13}CMS$ (Huang et al, 2014; Llufrio
64 et al, 2019) and PAVE (Wang et al, 2019). It is highly desirable to have a generic and flexible tool
65 to process untargeted isotope tracing metabolomics, and to enable more flexible data analysis
66 and modeling.

67

68 In this study, we propose a generalized tree structure to assign relationship of each ion to the
69 original compound and infer its neutral mass. The pre-annotation software tool, khipu, is freely
70 available as a Python package. It is applicable to both regular untargeted metabolomics and
71 stable isotope tracing data, and helps plug metabolomics data easily into common data science
72 tools.

73

74 **Results**

75 **The combination of isotopes and adducts is a 2-tier tree.**

76 The redundant or degenerate ions in mass spectrometry can be from in-source modifications
77 (adducts, fragments and conjugates) on any of the isotopic forms. For simplicity, we only consider
78 adducts in the initial steps. The combination of isotopes and adducts leads to a grid of mass
79 values, relative to the neutral mass of M_0 , exemplified in **Table 1**. We use M_0 to denote the
80 molecules with only ^{12}C atoms. The isotopes are denoted as $^{13}\text{C}/^{12}\text{C}$, $^{13}\text{C}/^{12}\text{C}^2$, etc., whereas
81 the last digit is the number of ^{13}C atoms present in each molecule.

82

83 The adducts can be represented as a tree (**Figure 1A**), using the neutral form as the root, which
84 is usually not measured in mass spectrometry. Each edge in the tree corresponds to a specific
85 mass difference, from the reaction forming the adduct. In fact, the full grid in Table 1 can be
86 accommodated into the tree, using isotopes as leaves to the adducts. Two arguments favor the
87 tree as a preferred data structure over a generic network: 1) each ion measured in mass
88 spectrometry is formed from a specific “predecessor”, and 2) the whole group of ions are from a
89 unique compound, which is the “root”. In computational terms, a network becomes a tree once it
90 fulfills the two requirements: 1) each node can have no more than one predecessor and 2) a
91 unique root. The benefit of this tree representation is important, allowing automated interpretation
92 of all ions via defined semantics.

93

94 Because the isotopes are present independently from each other at the time of measurement, we
95 treat them equally as one tier of the tree here. It is noted that the generation of them may have
96 biochemical significances in isotope tracing experiments, but that problem is outside data
97 processing and annotation. Therefore, the combination of isotopes and adducts, as exemplified
98 in Table 1, can be represented as a 2-tier tree. The tree can either use adducts as tier 1 or isotopes
99 as tier 1. The decision is to use adducts as tier 1, because a) adduct mass patterns are more
100 distinct, and b) isotopes are often limited by abundance, resulting only M_0 ions in many
101 compounds.

102

103 **An algorithm to convert a mass distance network to a 2-tier tree.**

104 Annotation methods in MS metabolomics commonly start by searching mass difference patterns,
105 e.g. 1.0034 for $^{13}\text{C}/^{12}\text{C}$ in isotopes and 22.9893 for Na^+ in adducts. Each match leads to a pair
106 of ions (also called features), and many pairs are connected via shared ions into a network of ions
107 (**Figure 1B**). During the mass difference search, additional redundancy is introduced, e.g. the
108 mass difference between ^{13}C and ^{12}C is the same as between $^{13}\text{C}/^{12}\text{C} \times 2$ and $^{13}\text{C}/^{12}\text{C} \times 3$, and
109 so forth. This network redundancy is apparent in the top part of the network in Figure 1B. The
110 objective in annotation is to identify the true root (original compound) from the network, which has
111 been challenging in previous works.

112

113 As biological reactions are not part of data annotation here, the edges in our mass distance
114 networks belong to one of the two categories: isotopic differences or in-source modifications
115 (**Figure 1B**). A key observation is that all ions connected by isotopic edges belong to the same
116 adduct. Therefore, subnetworks per adduct can be defined from a mass distance network (**Figure**
117 **1C**). Once these isotopic subnetworks are abstracted into individual network nodes, we can find
118 the best alignment between this abstracted network (**Figure 1C**) and the adduct tree (**Figure 1A**).
119 The algorithm is designed as two-step optimization: to obtain a tree with optimal number of ions
120 explained in the alignment of adduct trees, then in the alignment of isotopes. The result of this
121 algorithm on our example network is shown in **Figure 1D**. To match our 2-tier tree structure, the
122 networks have to become directed acyclic graph (DAG). During this process, erroneous edges
123 are weeded out because they do not satisfy DAG and a rooted tree. This method yields a
124 structured and unique annotation of each ion in the tree. Based on the matched m/z values, the
125 neutral mass of M_0 compound is obtained by a regression model. Once the core structure of a
126 tree is established, additional adducts and fragments can be searched in the data. The algorithm
127 is implemented into a freely available Python package *kipu*.

128

129 **Khipu plots allow intuitive interpretation of isotope tracing data.**

130 After ions are grouped into a tree for each original compound, they are recorded into transparent
131 JSON format, as defined for empirical compounds (see examples in **Supplemental notebook**).
132 An “empirical compound” refers to a tentatively defined compound in metabolomics data, used in
133 our previous projects (Li et al, 2013, Pang et al, 2020), as the technology may not deliver definitive
134 identification or resolve a mixture (e.g. isomers not successfully separated).

135

136 We continue using the compound in **Figure 1 B-D** to illustrate the khipu plotting functions. Each
137 ion is measured with an intensity value in one of more biological samples. While the tree
138 visualization in Figure 1D is useful, khipu includes multiple functions to visualize the features, *m/z*
139 values, intensity values as data frame tables (Supplemental notebook), to facilitate intuitive
140 interpretation of each compound. An enhance visualization of the tree is demonstrated in **Figure**
141 **2A**, where the adducts are organized as a “trunk” and isotopes as “branches”. It’s clear that
142 several isotopes are present as the protonated ion; Na and K adducts are present for the more
143 abundant isotopes. Because this visualization style resembles the khipu knot records used by
144 Andean South Americans, we named our software “khipu”.

145
146 This experimental dataset was from cultured *E. coli*, containing three unlabeled samples and three
147 samples grown on U-13C-glucose. **Figure 2B** visualizes the intensity values across samples for
148 the M+H⁺ ion. The three unlabeled samples have high M0 peaks, and smaller 13C/12C (M1)
149 peaks due to the naturally occurring isotopes. The U-13C labelled samples have the highest
150 peaks at 13C/12C*9 (M9), with smaller peaks of other isotopes. This indicates that the latter
151 samples are almost fully labelled by 13C, and the compound should contain 9 carbon atoms. The
152 neutral mass inferred by khipu is 187.1686, which matches to acetylspermidine, which has a
153 chemical formula C9H21N3O, perfectly consistent with the isotopic pattern.

154
155 As a pre-annotation tool, khipu is positioned to feed organized data for downstream data
156 analysis. Users can choose to model the isotopes and compute flux using other tools (Moseley
157 2010; Millard et al, 2012). Khipu results can be easily used by other software and analyzed
158 using common data science tools (demo notebooks included in the code repository). JSON
159 (JavaScript Object Notation) is a common format for data exchange between software programs
160 and web applications, and one of khipu export formats. This enables an effective way for
161 sharing metabolite annotation, which is human friendly, computable, and neutral to software
162 platforms. A snippet of khipu export in JSON is as follows:

```
163     {'interim_id': 'root@187.1686',  
164       'neutral_formula_mass': 187.1686,  
165       'MS1_pseudo_Spectra': [  
166         {'id': 'F2353',  
167           'mz': 188.1759,  
168           'rttime': 20.57,  
169           'representative_intensity': 25299447.0,  
170           'isotope': 'M0',  
171           'modification': 'M+H+',  
172           'ion_relation': 'M0,M+H+'},  
173         {'id': 'F1741',  
174           'mz': 197.2061,
```

```
175         'runtime': 20.57,  
176         'representative_intensity': 16395781.0,  
177         'isotope': '13C/12C*9',  
178         'modification': 'M+H+',  
179         'ion_relation': '13C/12C*9,M+H+'},  
180         ....  
181     ],  
182     'MS2_Spectra': []}  
183
```

184 **How many metabolites do we measure?**

185 Proper pre-annotation is key to answer the question of how many metabolites/compounds are
186 measured in an experiment, which is a matter that has been debated for over a decade. Many
187 studies overestimated the coverage because the database search was inflated by
188 redundant/degenerate features/ions. Studies from the Patti and Rabinowitz labs used isotope
189 tracing techniques, and suggested the numbers are around 1,000~2,000 in *E. coli* and yeast
190 (Mathieu and Patti, 2017; Wang et al, 2018). Our khipu software now provides systematic and
191 fast pre-annotation on metabolomic datasets.

192 In our *E. coli* data (reverse phase ESI+), 3,602 LC-MS features were measured, and khipu
193 annotated 548 empirical compounds (trees) from 1,745 features. Among the 548 empirical
194 compounds, 445 have multiple isotopes (**Figure 3A**). The remaining 1,857 features are
195 singletons, i.e., not grouped with any other features. In two yeast datasets from Rabinowitz lab,
196 khipu annotation resulted in 1,775 and 908 empirical compounds, respectively in ESI+ and ESI-
197 modes (**Figure 3B&C**). In the yeast datasets, we included additional adducts from Wang et al
198 (2021), which by design did not increase the number of empirical compounds, but increased the
199 explained ESI+ features from 6,310 to 8,049, and from ESI- features 2,601 to 2,912. These results
200 suggest that less than 2,000 compounds were reliably measured in these experiments. Of note,
201 closer examination of each dataset should also remove contaminants, which is not part of khipu.

202

203 **Discussion**

204 Annotation of untargeted metabolomics data, including isotopic tracing data, is still not fully
205 solved. Many current tools take a network approach but depend on assumed base ions or
206 formulas to assign relationship between ions. We present a new algorithm here to resolve the
207 mass distance networks into a tree structure, unambiguously defining ion relationships and
208 inferring neutral mass. This approach shall reduce false annotations, and facilitate new compound
209 identification and discoveries. We consider the pre-annotation with khipu a key step forward, also
210 because it ships with generalized annotation format, which will greatly facilitate data exchange

211 and software interoperability. With this foundation, future benchmarking and improvements are
212 expected.

213
214 Multiple Jupyter notebooks are provided as part of the software package to demonstrate how
215 khipu is plugged into common data science tools. This gives great flexibility to people in using
216 both regular and isotope tracing metabolomics data, because the computational methods, as well
217 as experimental designs, are no longer limited by rigid software designs. This is an emerging
218 model in data science. Traditional software development is often too costly, and its maintenance
219 is even more challenging (Chang et al, 2021). Fundamentally, no software developer can meet
220 every demand via point-and-click interface. Scientific data analysis has to depend greatly on
221 scripting. The combination of modular software components, transparent data structures and
222 Jupyter notebooks opens up many opportunities for collaborations and scientific progress (Pittard
223 et al, 2020).

224
225 Khipu can be easily reused by other software tools. We plan to integrate it with the preprocessing
226 software asari (Li et al, 2022), whereas elution patterns can be better determined than from
227 standalone feature tables. The standard input to khipu is tab delimited feature tables, which
228 should be compatible with any LC-MS preprocessing software. Therefore, it will be easy to
229 incorporate it into metabolomics workflows, where complete annotation can take into
230 consideration of contaminants, authentic libraries and tandem mass spectrometry data.

231

232

233 **Methods**

234 **Python implementation:** Khipu is developed as an open source Python 3 package, and available
235 to install from the standard PyPi repository via the pip tool. It is freely available on GitHub
236 (<https://github.com/shuzhao-li/khipu>) under a BSD 3-Clause License. The graph operations are
237 supported by the networkx library, tree visualization aided by the treelib library. Khipu uses our
238 package mass2chem for search functions. The data model of “empirical compound” is described
239 in the metDataModel package. The package is designed in a modular way to encourage reuse.
240 The classes of Weavor and Khipu contain main algorithms, supported by numerous utility
241 functions. All functions are documented in the source via docstrings. Examples of reuse are given
242 in wrapper functions and in Jupyter notebooks. It can be run as a standalone command line tool.
243 Users can use a feature table from any preprocessing tool as input and get annotated empirical
244 compounds in JSON and tab delimited formats.

245

246 **LC-MS metabolomics data.** The dry extracts of unlabeled and ^{13}C labeled *E. coli* (Cambridge
247 Isotope Laboratories, Inc.; Catalog number: MSK-CRED-DD-KIT) were reconstituted in 100 μL of
248 $\text{ACN}/\text{H}_2\text{O}$ (1:1, v/v) then sonicated (10 mins) and centrifuged (10 mins at 13,000 rpm and 4°C)
249 before overnight incubation at 4°C . The supernatant for each $^{12}\text{C}/^{13}\text{C}$ *E. coli* extract was collected
250 and then prepared for LC-MS analysis. Metabolite extraction was carried out using
251 acetonitrile:methanol (8:1, v/v) containing 0.1% formic acid. All samples were vortexed and
252 incubated with shaking at 1000 rpm for 10 min at 4°C followed by centrifugation at 4°C for 15 min
253 at 15,000 rpm. The supernatant was transferred into mass spec vials and 2 μl injected into
254 UHPLC-MS. All samples were maintained at 4°C in the autosampler, and analyzed using a
255 Thermo Scientific Orbitrap ID-X Tribrid Mass Spectrometer coupled to a Thermo Scientific
256 Transcen LX-2 Duo UHPLC system, with a HESI ionization source, using positive ionization. A
257 Hypersil GOLDTM RP column (3 μm , 2.1 mm x 50 mm) maintained at 45°C was used. 0.1%
258 formic acid in water and 0.1% formic acid in acetonitrile were used as mobile phase A and B
259 respectively. The following gradient was applied at a flow rate of 0.4 ml/min: 0-0.1 min: 0% B,
260 0.10-1.9 min: 60% B, 1.9-5.0 min: 98% B, 5.00-5.10 min: 0% B and 4.9 min cleaning and column
261 equilibration. The chromatographic run time was 5 min followed by 5 min washing step after each
262 sample. The MS settings are: spray voltage, 3500 V; sheath gas, 45 Arb; auxiliary gas, 20 Arb;
263 sweep gas, 1 Arb; ion transfer tube temperature, 325°C ; vaporizer temperature, 325°C ; mass
264 range, 80-1000 Da; maximum injection time, 100 ms; resolution 60,000.

265

266 The yeast data from Chen et al. 2021 were retrieved from the Massive repository
267 (<https://massive.ucsd.edu>, ID no. MSV000087434). The yeast ESI+ data contain both unlabeled
268 and ^{13}C isotope labeled samples, while the ESI- data did not involve isotope tracing, The data
269 from Mathieu and Patti (2017) and Wang et al (2018) were not found publicly. All datasets were
270 processed using asari version 1.9.2 (<https://github.com/shuzhao-li/asari>). The yeast ESI- dataset
271 was quality filtered for signal-noise-ratio > 100 to serve as a cleaner demo.

272

273 **Tables**

274

275 **Table 1. Combinations of isotopes and adducts generate mass differences as a grid.**

276 The mass values are relative to the ^{12}C only neutral mass. Examples are using a limited

277 number of isotopes and in-source modifications in positive ionization.

278

	<i>M+H[+]</i>	<i>M+NH4[+]</i>	<i>M+Na[+]</i>	<i>M+HCl+H[+]</i>	<i>M+K[+]</i>	<i>M+ACN+H[+]</i>
<i>MO</i>	1.007276	18.033826	22.989276	36.983976	38.963158	42.033825
<i>$^{13}\text{C}/^{12}\text{C}$</i>	2.010631	19.037181	23.992631	37.987331	39.966513	43.03718
<i>$^{13}\text{C}/^{12}\text{C}^*2$</i>	3.013986	20.040536	24.995986	38.990686	40.969868	44.040535
<i>$^{13}\text{C}/^{12}\text{C}^*3$</i>	4.017341	21.043891	25.999341	39.994041	41.973223	45.04389
<i>$^{13}\text{C}/^{12}\text{C}^*4$</i>	5.020696	22.047246	27.002696	40.997396	42.976578	46.047245
<i>$^{13}\text{C}/^{12}\text{C}^*5$</i>	6.024051	23.050601	28.006051	42.000751	43.979933	47.0506
<i>$^{13}\text{C}/^{12}\text{C}^*6$</i>	7.027406	24.053956	29.009406	43.004106	44.983288	48.053955

279

280 **Figures**

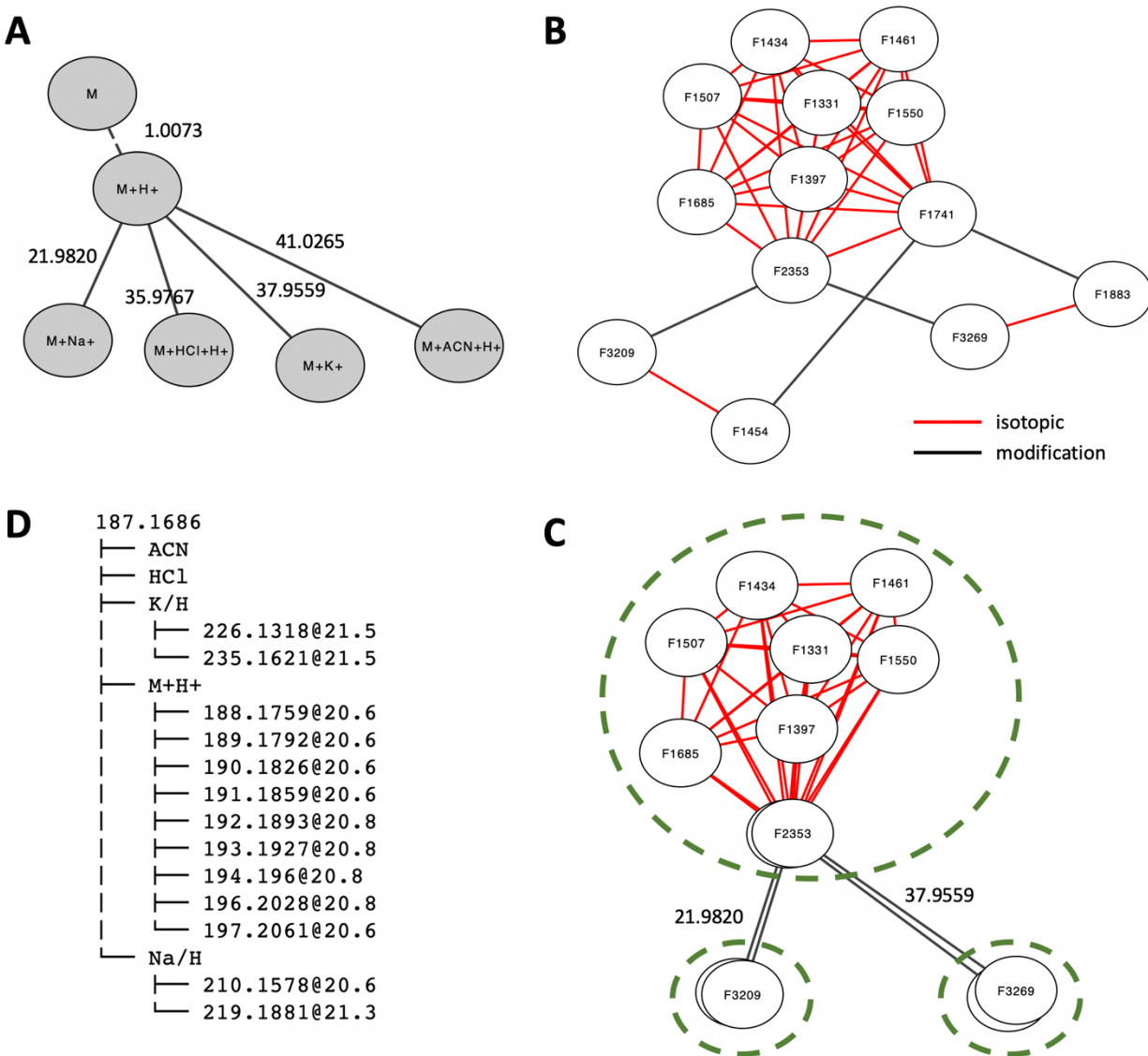
281 **Figure 1. The khipu algorithm converts a mass distance network to a tree structure.**

282 A) An adduct tree base on Table 1. Mass differences on the edges are relative to the
283 predecessor nodes.

284 B) An example mass distance network from our credentialed E. coli dataset, which contains
285 both unlabeled and ¹³C labeled samples. Edges in red are from isotopic patterns and edges in black
286 from adduct patterns.

287 C) The isotopic subnetworks can be treated as individual nodes, then the abstracted network
288 has only adduct edges, which facilitates the alignment to the theoretical adduct tree in A).

289 D) Resulted 2-tier tree. The root is inferred neutral mass. No ion is assigned to ACN or HCl
290 adducts. Decimal numbers should be consistent with that in Table 1.



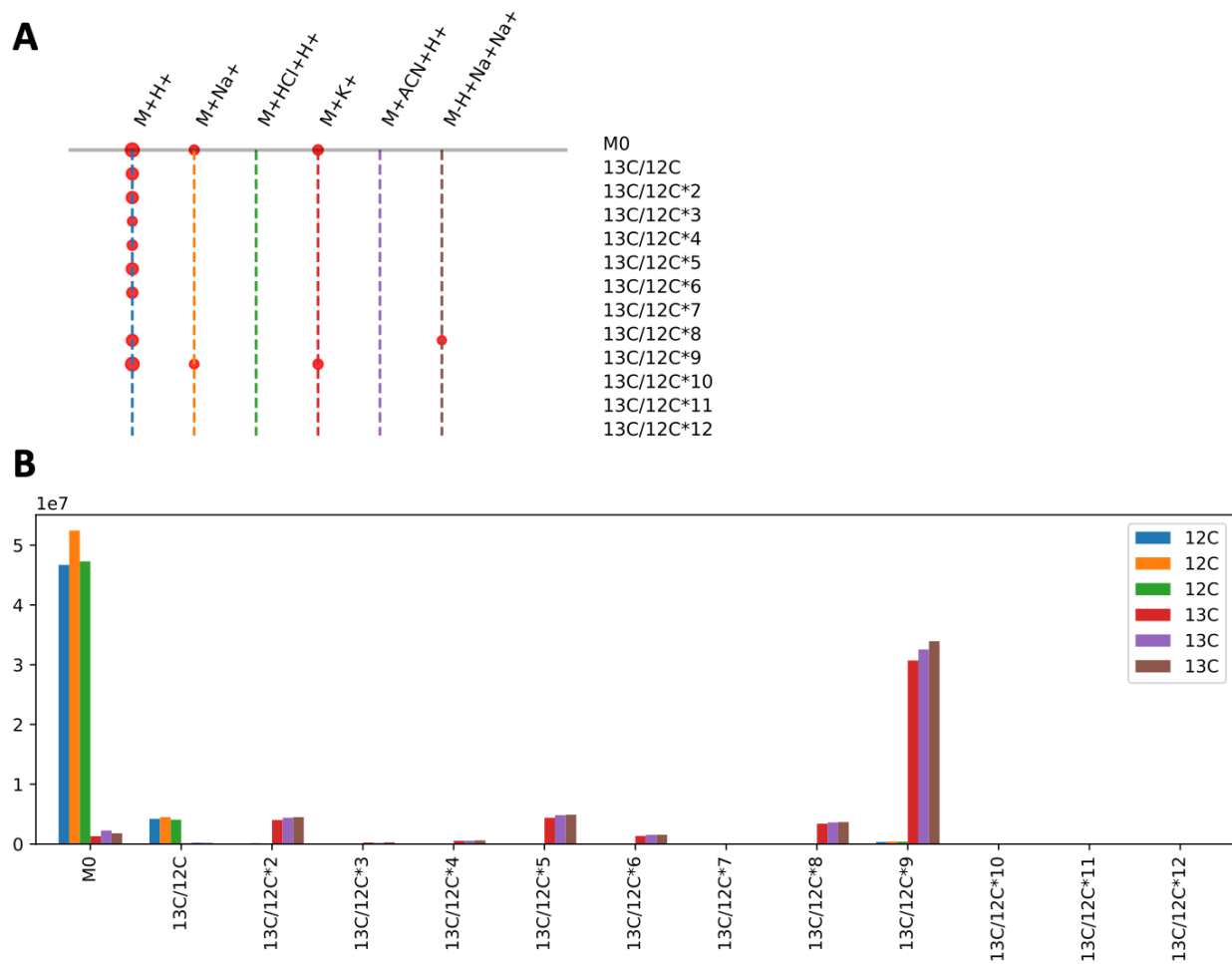
292

293 **Figure 2. Visualization using khipu facilitates interpretation of isotope tracing data.**

294 A) An example khipugram plot for the compound in Figure 1, with its 13 ions aligned to the tree
295 in Figure 1D. Each dot represents an ion measured in the data, the size of dots proportional to
296 average intensity. The vertical dashed lines are colored for easy navigation, and the colors are
297 of no particular meaning.

298 B) Bar plot for intensity values of the $M+H^+$ ion in different isotopes (x-axis) for three ^{12}C
299 samples and three ^{13}C samples (in color legend). This is from the first branch in A).

300



301

302

303

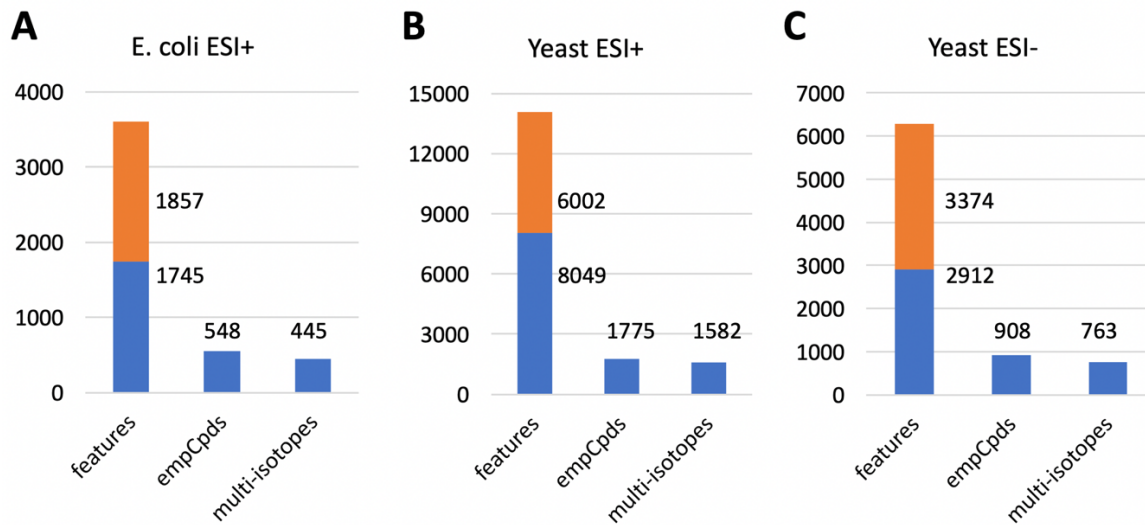
304

305

306

307

308 **Figure 3. Number of measured compounds in three metabolomic datasets.**
309 A) Credentialed *E. coli* data generated in this study. B) Previously published yeast ESI+ and C)
310 ESI- datasets from Rabinowitz lab (Wang et al. 2021). Khipu annotation on these datasets took
311 2~6 seconds on a laptop computer of Intel i7 CPU. The orange portions are referred as
312 “singletons”.
313



314
315

316 **Supplemental File:**

317

318 **data_analysis_ecoli_pos.pdf**

319 A Jupyter Notebook printed to PDF format to demonstrate khipu applications.

320

321 **References:**

- 322
- 323 Blaženović, I., Kind, T., Sa, M.R., Ji, J., Vaniya, A., Wancewicz, B., Roberts, B.S., Torbašinović,
324 H., Lee, T., Mehta, S.S. and Showalter, M.R., 2019. Structure annotation of all mass
325 spectra in untargeted metabolomics. *Analytical chemistry*, 91(3), pp.2155-2162.
326
- 327 Bueschl, C., Kluger, B., Neumann, N.K., Doppler, M., Maschietto, V., Thallinger, G.G., Meng-
328 Reiterer, J., Krska, R. and Schuhmacher, R., 2017. MetExtract II: a software suite for
329 stable isotope-assisted untargeted metabolomics. *Analytical chemistry*, 89(17), pp.9518-
330 9526.
331
- 332 Chang, H.Y., Colby, S.M., Du, X., Gomez, J.D., Helf, M.J., Kechris, K., Kirkpatrick, C.R., Li, S.,
333 Patti, G.J., Renslow, R.S. and Subramaniam, S., 2021. A practical guide to
334 metabolomics software development. *Analytical chemistry*, 93(4), pp.1912-1923.
335
- 336 Chen, L., Lu, W., Wang, L., Xing, X., Chen, Z., Teng, X., Zeng, X., Muscarella, A.D., Shen, Y.,
337 Cowan, A. and McReynolds, M.R., 2021. Metabolite discovery through global annotation
338 of untargeted metabolomics data. *Nature methods*, 18(11), pp.1377-1385.
339
- 340 Chokkathukalam, A., Jankevics, A., Creek, D.J., Achcar, F., Barrett, M.P. and Breitling, R.,
341 2013. mzMatch–ISO: an R tool for the annotation and relative quantification of isotope-
342 labelled mass spectrometry data. *Bioinformatics*, 29(2), pp.281-283.
343
- 344 DeFelice, B.C., Mehta, S.S., Samra, S., Cajka, T., Wancewicz, B., Fahrman, J.F. and Fiehn,
345 O., 2017. Mass spectral feature list optimizer (MS-FLO): a tool to minimize false positive
346 peak reports in untargeted liquid chromatography–mass spectroscopy (LC-MS) data
347 processing. *Analytical chemistry*, 89(6), pp.3250-3255.
348
- 349 Domingo-Almenara, X., Montenegro-Burke, J.R., Benton, H.P. and Siuzdak, G., 2018.
350 Annotation: a computational solution for streamlining metabolomics analysis. *Analytical*
351 *chemistry*, 90(1), p.480.
352
- 353 Huang, X., Chen, Y.J., Cho, K., Nikolskiy, I., Crawford, P.A. and Patti, G.J., 2014. X13CMS:
354 global tracking of isotopic labels in untargeted metabolomics. *Analytical chemistry*, 86(3),
355 pp.1632-1639.
356
- 357 Kachman, M., Habra, H., Duren, W., Wigginton, J., Sajjakulnukit, P., Michailidis, G., Burant, C.
358 and Karnovsky, A., 2020. Deep annotation of untargeted LC-MS metabolomics data with
359 Binner. *Bioinformatics*, 36(6), pp.1801-1806.
360
- 361 Kuhl, C., Tautenhahn, R., Bottcher, C., Larson, T.R. and Neumann, S., 2012. CAMERA: an
362 integrated strategy for compound spectra extraction and annotation of liquid
363 chromatography/mass spectrometry data sets. *Analytical chemistry*, 84(1), pp.283-289.
364
- 365 Li, Shuzhao, Amnah Siddiq, Maheshwor Thapa, and Shujian Zheng. "Trackable and scalable
366 LC-MS metabolomics data processing using asari." *bioRxiv* (2022). doi:
367 <https://doi.org/10.1101/2022.06.10.495665>
368
- 369 Li, S., Park, Y., Duraisingham, S., Strobel, F.H., Khan, N., Soltow, Q.A., Jones, D.P. and
370 Pulendran, B., 2013. Predicting network activity from high throughput metabolomics.
371 *PLoS computational biology*, 9(7), p.e1003123.

- 372
373 Llufrío, E.M., Cho, K. and Patti, G.J., 2019. Systems-level analysis of isotopic labeling in
374 untargeted metabolomic data by X13CMS. *Nature protocols*, 14(7), pp.1970-1990.
375
376 Mahieu, N.G. and Patti, G.J., 2017. Systems-level annotation of a metabolomics data set
377 reduces 25 000 features to fewer than 1000 unique metabolites. *Analytical chemistry*,
378 89(19), pp.10397-10406.
379
380 Mahieu, N.G., Spalding, J.L., Gelman, S.J. and Patti, G.J., 2016. Defining and detecting
381 complex peak relationships in mass spectral data: the Mz. unity algorithm. *Analytical*
382 *chemistry*, 88(18), pp.9037-9046.
383
384 Millard, P., Letisse, F., Sokol, S. and Portais, J.C., 2012. IsoCor: correcting MS data in isotope
385 labeling experiments. *Bioinformatics*, 28(9), pp.1294-1296.
386
387 Moseley, H.N., 2010. Correcting for the effects of natural abundance in stable isotope resolved
388 metabolomics experiments involving ultra-high resolution mass spectrometry. *BMC*
389 *bioinformatics*, 11(1), pp.1-6.
390
391 Naake, T. and Fernie, A.R., 2018. MetNet: metabolite network prediction from high-resolution
392 mass spectrometry data in R aiding metabolite annotation. *Analytical chemistry*, 91(3),
393 pp.1768-1772.
394
395 Pang, Z., Chong, J., Li, S. and Xia, J., 2020. MetaboAnalystR 3.0: toward an optimized workflow
396 for global metabolomics. *Metabolites*, 10(5), p.186.
397
398 Pittard, W.S., Villaveces, C. and Li, S., 2020. A Bioinformatics Primer to Data Science, with
399 Examples for Metabolomics. *Computational Methods and Data Analysis for*
400 *Metabolomics*, pp.245-263.
401
402 Previs, S.F. and Downes, D.P., 2020. Key Concepts Surrounding Studies of Stable Isotope-
403 Resolved Metabolomics. In *Computational Methods and Data Analysis for Metabolomics*
404 (pp. 99-120). Humana, New York, NY.
405
406 Rahim, M., Ragavan, M., Deja, S., Merritt, M.E., Burgess, S.C. and Young, J.D., 2022. INCA
407 2.0: A tool for integrated, dynamic modeling of NMR-and MS-based isotopomer
408 measurements and rigorous metabolic flux analysis. *Metabolic Engineering*, 69, pp.275-
409 285.
410
411 Senan, O., Aguilar-Mogas, A., Navarro, M., Capellades, J., Noon, L., Burks, D., Yanes, O.,
412 Guimera, R. and Sales-Pardo, M., 2019. CliqueMS: a computational tool for annotating
413 in-source metabolite ions from LC-MS untargeted metabolomics data based on a
414 coelution similarity network. *Bioinformatics*, 35(20), pp.4089-4097.
415
416 Uppal, K., Walker, D.I. and Jones, D.P., 2017. xMSannotator: an R package for network-based
417 annotation of high-resolution metabolomics data. *Analytical chemistry*, 89(2), pp.1063-
418 1067.
419
420 Wang, L., Xing, X., Chen, L., Yang, L., Su, X., Rabitz, H., Lu, W. and Rabinowitz, J.D., 2018.
421 Peak annotation and verification engine for untargeted LC-MS metabolomics. *Analytical*
422 *chemistry*, 91(3), pp.1838-1846.

423

424 **Code availability:** The asari source code is available at GitHub, [https://github.com/shuzhao-](https://github.com/shuzhao-li/khipu)
425 [li/khipu](https://github.com/shuzhao-li/khipu), and as a Python package via <https://pypi.org/project/khipu-metabolomics/>. The
426 demonstration datasets are provided as part of the source code.

427

428 **Acknowledgements:** This work was in part funded by NIH grants (to SL) U01 CA235493
429 (NCI) and R01 AI149746 (NIAID).

430

431 **Author contributions:** S.L. designed the study, wrote the khipu software and the manuscript.
432 S.Z. performed the LC-MS metabolomics experiment on credentialed E. coli samples.

433