

1
2 Editorial summary: MASST+ speeds up querying of metabolomics mass spectrometry data by two
3 orders of magnitude.

4

5 Fast Mass Spectrometry Search and Clustering of Untargeted 6 Metabolomics Data

7
8 Mihir Mongia^{*1}, Tyler M. Yasaka^{*1}, Yudong Liu^{*1}, Mustafa Guler¹, Liang Lu¹, Aditya Bhagwat¹,
9 Bahar Behsaz^{1,2}, Mingxun Wang³, Pieter C. Dorrestein^{4,5}, and Hosein Mohimani¹

10
11¹Computational Biology Department, School of Computer

12 Science Carnegie Mellon University

13
14²Chemia Biosciences Inc.

15
16³Computer Science and Engineering, University of California Riverside

17
18⁴Collaborative Mass Spectrometry Innovation Center, Skaggs School of Pharmacy
19 and Pharmaceutical Sciences, University of California San Diego

20
21⁵Department of Pharmacology and Pediatrics, University of California San Diego

22

23 The throughput of mass spectrometers and the size of publicly available metabolomics data are
24 growing rapidly, but analysis tools like molecular networking and MASST do not scale to searching
25 and clustering billions of mass spectral data in metabolomics repositories. To address this
26 limitation, we designed MASST+ and Networking+, which can process datasets that are up to three
27 orders of magnitude larger than the state-of-the-art.

28

29

30

Introduction

31 During the past decade, the size of mass spectral data collected in the fields of natural products,
32 exposomics, and metabolomics has grown exponentially^{9,16,18}. In accordance with the advances in
33 mass spectrometry technology, multiple computational methods were developed for analyzing this
34 massive data. Recently Mass Spectrometry Search Tool (MASST) was introduced as a search
35 engine for finding analogs of a query spectrum in mass spectrometry repositories¹⁹. MASST has
36 demonstrated utility in the annotation of a wide variety of unidentified metabolites, including
37 clinically important molecules in patient cohorts^{15,3,6}, toxins/pesticides in environmental samples¹⁴,
38 fungal metabolites¹⁰, and metabolites from pathogenic microorganisms^{4,11,5}. Moreover, molecular
39 networking was introduced for clustering spectral datasets into families of related molecules^{28,29}.
40 Molecular Networking has yielded a systematic view of the chemical space in different ecosystems
41 and helped determine the structure of many compounds^{20,21,22,23,24,25,27,26}.

42 MASST and molecular networking are based on a naive approach for scoring two tandem mass
43 spectra. MASST compares the query spectrum against all reference spectra one by one and
44 computes a similarity score based on the relative intensities of shared and shifted peaks. Therefore,
45 the runtime of MASST grows linearly with the repository size. Molecular networking first uses MS-
46 Clustering²⁸ to cluster identical spectra by calculating a dot-product score (ExactScore, Figure 1a-i)
47 between the spectra. Then Spectral Networking²⁹ is used to calculate a dot product score that
48 accounts for peaks that are shared or shifted (ShiftedScore, Figure 1a-ii) between all pairs of
49 clusters in order to find groups of related molecules. This latter procedure grows quadratically with
50 the number of clusters. Current trends show that the size of public mass spectral repositories
51 doubles every two to three years (Supplementary Fig. 1). Therefore, the current implementations of
52 MASST and Molecular Networking will not be able to scale with the growth of future repositories.
53 A MASST search for a single spectrum against the clustered global natural product social (GNPS)
54 database (~83 million clusters) currently takes about an hour on a single thread and a MASST

55 search against the entire GNPS (717 million spectra) does not complete after being run for three
56 days. Currently, molecular networking analysis of a million spectra takes a few hours, while
57 molecular networking of ~20 million spectra does not yield results after running for a week. Similar
58 to the area of computational genomics, handling the exponential growth of repositories requires the
59 development of more efficient and scalable search algorithms.

60 In this paper, we introduce a fast dot product algorithm that preprocesses a set of spectra into an
61 indexing table. This indexing table maps all possible precursor m/z and fragment ion m/z pairs to
62 the spectra that contain them. Using this indexing, given a query spectrum, the dot product with
63 respect to all spectra can be computed efficiently by iterating through each query peak and using the
64 indexing table to retrieve spectra with similar peaks (Figure 1b). Since mass spectra are sparse, only
65 a small fraction of spectra/peaks are retrieved for each query. The ability to leverage this sparsity
66 requires only a small fraction of the compute used by naive scoring methods because the vast
67 majority of the MS/MS spectra in the index are never touched during the query process. By
68 integrating this indexing approach into the scoring subroutines of MASST and Molecular
69 Networking, we develop two computational tools, MASST+ and Networking+, that are two to three
70 orders of magnitude faster than state-of-the-art on large datasets. Further, the indexing approach
71 supports on-line growth, that is, the insertion of new spectra without the need for recalculation from
72 scratch. This enables both MASST+ and Networking+ to efficiently handle the dynamic growth of
73 reference spectra. Currently MASST+ is available as a web service from
74 <https://masst.ucsd.edu/masstplus/>. GNPS supports stand-alone MASST+ (Supplementary Fig. 2)
75 and integration with molecular networking (Supplementary Fig. 3).

76

77

78

79

80

81

82

Results.

83 **Outline of MASST+ algorithm.** Given a query spectrum, MASST+ efficiently searches a database
84 of reference spectra to find similar entries by creation of an indexing table – a data structure which
85 allows rapid retrieval of similar spectra based on the peaks present in the query spectrum. For each
86 precursor mass M and each peak mass p , a list of indices of spectra with precursor M and peak p are
87 stored, along with the intensity of the peaks. In case of exact search, MASST+ iterates through the
88 peaks in the query spectrum and retrieves the lists associated with a query peak and query's
89 precursor mass. The ExactScore is calculated by multiplying and adding up the intensity of each
90 peak in query spectrum and reference spectra (Figure 1b). In case of analog search (Supplementary
91 Fig. 4), MASST+ uses a much larger precursor mass tolerance (e. g. 300Da) and computes
92 ShiftedScore that takes into account both shared and Δ -shifted peaks (peaks in reference spectra that
93 are Δ Da larger than peaks in query), where Δ is the mass difference between the precursor of query
94 and reference spectra (Figure 1c).

95

96

97

98

99 **Outline of Networking+ algorithm.** Networking+ clusters spectral datasets into families of related
100 molecules by first putting spectra from identical molecules into the same clusters (Clustering+),
101 then forming the centers of each cluster by taking their consensus, and then connecting the clusters
102 that are predicted to be generated from related molecules (Pairing+). Clustering+ iterates over all
103 spectra and associates each spectrum with a cluster that is highly similar. It uses a strategy similar
104 to MASST+ exact search for efficiently calculating the SharedScore between the spectrum and each

105 cluster center. Pairing+ uses a shared and Δ -shifted dot-product as a similarity measure for
106 identifying related spectra. It uses a strategy similar to MASST+ analog search to find all pairs of
107 clusters with high ShiftedScore.

108 **Benchmarking MASST+.** We have benchmarked MASST+ (Supplementary Table 1) on various
109 GNPS datasets including MSV000078787 dataset collected on *Streptomyces* cultures (5,433
110 spectra), clustered GNPS (83,131,248 spectra), and entire GNPS (717,395,473 spectra).

111 Supplementary Data 1 lists Accession IDs of all GNPS datasets used in our study. While MASST
112 and MASST+ report identical hits, MASST+ is two orders of magnitude faster and more memory
113 efficient (Supplementary Table 1). For small data sets we only get a 3-fold increase in speed. This
114 becomes magnified when the data set that is searched becomes larger. In case of the clustered
115 GNPS, MASST+ performs analog search in 15 seconds while MASST takes 49 min, a 196-fold
116 increase. In case of the entire GNPS, MASST+ performs analog search in under two hours on
117 average, while MASST search does not finish within three days on the GNPS server making it
118 practically not possible to routinely perform such a search.

119 Figure 2a illustrates the runtime and memory consumption of MASST+ in exact and analog mode
120 for various subsets of the clustered GNPS. Supplementary Fig. 5 illustrates that indexing time and
121 memory consumption grows linearly with the size of datasets and Supplementary Fig. 6 shows
122 indexing time increases for larger values of peak mass tolerance. MASST+ takes eight hours of
123 compute time and eight gigabytes of memory to index ~83 million spectra from the clustered GNPS
124 and 72 hours of compute time and 9 gigabytes of memory to index 717 million spectra contained in
125 GNPS. Supplementary Fig. 7 breaks down MASST+ runtime into two different steps, loading
126 peaks lists and computing dot product, for various numbers of query spectra. Loading peak lists
127 consumes about half of the total runtime when the number of query spectra is greater than 100.

128

129

130
131
132
133
134
135
136
137
138
139
140
141
142

143

144 **Benchmarking networking+**. Figure 2b, Supplementary Table 2 and Supplementary Tables 3-5
145 benchmark Networking+ against molecular networking on various data sizes for which runtime is
146 less than 24 hours. In 24 hours Clustering+ can process 300 million spectra on a single CPU, while
147 MS-Clustering can process 20 million spectra. Moreover, in this timeline, Pairing+ can process 2
148 million spectra, while spectral networking can handle 0.2 million spectra. Clustering+ and Pairing+
149 are two orders of magnitude faster than their counterparts, MS-Clustering²⁹ and Spectral
150 Networking²⁸. The clusters and networks reported by Clustering+ and Pairing+ are identical to MS-
151 Clustering and spectral networks. As previously noted in Bittremieux et al.⁴³, it was not possible to
152 directly create a molecular network from all the GNPS spectra, here we show that this is now
153 possible with Networking+ with minimal computer memory requirements.

154

155
156
157
158
159

160

161 **Networking the entire GNPS.** We clustered the entire GNPS (717 million scans) using
162 Clustering+ and formed the network using Pairing+. This resulted in 8,453,822 million clusters and
163 4,947,928 connected components with a total of 17,533,386 edges (available from

164 https://github.com/mohimanilab/MASSTplus). Among 4,948,146 connected components in the
165 network, 98% (4,849,047 components) consist of a single node, while 1.5%, 0.3%, 0.2% and 0.02%
166 (74530, 13957, 9239, and 1152 components) had 2, 3, and 4-9 and 10+ nodes (Supplementary Fig.
167 8). Among 7,986,356 clusters in the network, 1.7% (134,198 Clusters) matched reference spectra
168 from the NIST library, 6% (477,721 clusters) were a neighbor of a cluster matched NIST library,
169 14% (1,130,092 clusters) were a neighbor of a neighbor, and 78% (5,390,554 clusters) were three or
170 more hops away from any cluster matching NIST library (Supplementary Fig. 9). Of 307,709
171 clusters consisting of 20 or more spectra, for 18% (54,518 clusters) all spectra came from a single
172 MassIVE dataset, while for 13% and 69% (39,428 and 213,763 clusters) spectra came from 2 or 3+
173 MassIVE datasets (Supplementary Fig. 10). About 61 percent of the clusters with precursor mass
174 between 0 and 400 Daltons consisted of only two GNPS spectra whereas less than half the clusters
175 with precursor mass above 400 Daltons consisted of only two GNPS spectra (Supplementary Fig.
176 11). Networking+ took 6 days to finish this task on 1 CPU. Currently, this task is not feasible using
177 existing approaches.

178

179

180 **Applying Networking+ for Identification of lanthipeptides.** The indexing strategies proposed
181 here are applicable to all classes of small molecules. Here we illustrate the application of these
182 methods in the case of lanthipeptide natural products. Currently, methods for high-throughput
183 discovery of lanthipeptides through computational analysis of genomics and metabolomics data
184 suffer from various limitations, especially at repository scale. Lanthipeptides are a biologically
185 important class of natural products that include antibiotics³⁰, antifungals³¹, antivirals³², and
186 antinociceptives³³. Lanthipeptides are structurally defined by the thioether amino acids lanthionine,
187 methyllanthionine and labionin. Lanthionine and methyllanthionine are introduced by dehydration
188 of a serine or threonine (to generate a dehydroalanine or dehydrobutyrine) and addition of a
189 cysteine thiol, catalyzed by a dehydratase and a cyclase, respectively³⁴. During lanthipeptide

190 biosynthesis, a precursor gene lanA is translated by the ribosome to yield a precursor peptide LanA
191 that consists of a N-terminal leader peptide and a C-terminal core peptide sequence. The core
192 peptide is post-translationally modified by the lanthionine biosynthetic machinery and other
193 enzymes. It is then proteolytically cleaved from the leader peptide to yield the mature lanthipeptide
194 and exported out of the cell by transporters.

195 Lanthipeptides usually possess network motifs that enable mining them in spectral networks. These
196 motifs include mass shifts of -18.01Da (H₂O mass) that correspond to the varying number of
197 dehydrations, and mass shifts equal to amino acid masses that correspond to promiscuity in N-
198 terminal leader processing. We formed the spectral network using Networking+ for a subset of 500
199 *Streptomyces* cultures with known genomes (Supplementary Table 6). The dataset contains
200 9,410,802 scans, which are clustered into 354,401 nodes, 6,032 connected components, and
201 1,265,311 edges. Currently, Molecular Networking crashes on this dataset after eight days of
202 processing. We further only retained 29,639 nodes that possess the network motif by filtering for
203 edges with mass differences equal to a loss of H₂O, NH₃, or an amino acid mass. Then we filtered
204 for nodes with long amino acid sequence tags of various lengths using PepNovo³⁵ (Supplementary
205 Table 7). There are a total of 2,353 nodes with sequence tags of length 12 or longer, and 285 of
206 these nodes are connected to an edge with a mass difference equal to the mass of one H₂O or an
207 amino acid loss. We further inspected these nodes using our in-house software algorithm, Seq2Ripp
208 (<https://github.com/mohimanilab/seq2ripp>). Given a lanthipeptide precursor, Seq2Ripp generates all
209 molecular structures of all possible candidate molecules by considering different cores and various
210 modifications and then searches the candidate molecular structures against mass spectra using
211 Dereplicator³⁶. This strategy identified three known and 14 novel lanthipeptides with p-values
212 below 1e-15 (Supplementary Table 8). Among them, the precursor of 13 lanthipeptides (76%)
213 overlaps with reports by the genome mining strategy introduced by Walker et al.⁴⁰. However, only
214 for two lanthipeptides, the core peptides predicted are consistent with predictions from Walker et al.
215 (11%). Note that in contrast to our approach, Walker et al. is based solely on genomics, and it does

216 not use metabolomics data for identifying the start of core peptide. This demonstrates that MASST+
217 and Molecular Networking+ can be used to gain insight into previously uncharacterized molecules. One of
218 the novel peptides (CHM-1731 from *Streptomyces albus*) is further described in Figure 2c.

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238

239

240

241

242 **Discussion.**

243 The mass spectrometry search tool (MASST) and molecular networking have become powerful
244 strategies to analyze LC-MS/MS based data to a broad range of users in the research
245 community^{2,13,15,17,37,38,39}. However, these tools do not scale to searching and clustering large
246 spectral repositories with hundreds of millions of spectra. As the size of mass spectral repositories
247 doubles every two to three years, the current implementation of MASST and Molecular Networking
248 will soon not be able to meet the needs of biologists and clinicians and thus new solutions are
249 urgently needed.

Recent advances have enabled the determination of molecular formula⁴⁴ and chemical class⁴⁵ for a large portion of spectra in GNPS. Despite these efforts, it is challenging to assign a chemical structure to the majority of spectra in GNPS. MASST+ and Networking+ provide efficient ways to annotate this dark matter by elucidating known molecules and their novel variants in repositories as they grow to billions of mass spectra. MASST+ currently searches query spectra against the clustered GNPS in a few seconds (in comparison to an hour for MASST), hence enabling instant analysis of the query mass spectrum of interest. Further, MASST+ can search the entire GNPS, which contains hundreds of millions of spectra in less than two hours, a task that is currently impossible with MASST. MASST+ can be parallelized by splitting a set of query spectra among several computational nodes/threads. Each thread then can run a separate MASST+ search job that utilizes the same index stored on disk.

261

262 **Acknowledgements.**

263 The work of T.Y., M.M., Y.D., B.B., and H.M. was supported by a National Institutes of Health New
264 Innovator Award DP2GM137413, a U.S. Department of Energy award DE- SC0021340, a National
265 Science Foundation award DBI-2117640, and a National Institute of General Medicine Sciences of the
266 National Institutes of Health award R43GM150301 (B.B. only). The work of P.C.D and M.W. were
267 supported by R03OD034493, U24DK133658, and R01GM107550 (P.C.D only).

268

269 **Author Contributions Statement.**

270 M.M., T.M.Y., Y.L., M.G., L.L., and A.B. implemented the algorithms. M.M., T.M.Y., and Y.L. performed the
271 analysis. M.W. designed and implemented the GNPS web service for MASST+. B.B., P.C.D., and H.M.
272 designed and directed the work. M.M. and H.M. wrote the manuscript. All the co-authors contributed to the
273 revision of the manuscript.

274

275 **Competing Interests.**

276 H.M. and B.B. are co-founders and have equity interests from Chemia.ai, LLC. PCD is an advisor and
277 holds equity in Cybele, consulted for MSD animal health in 2023, and he is a Co-founder, holds equity
278 in and is scientific advisor for Ometa Labs, Arome, and Enveda with prior approval by UC San Diego.
279 The other authors declare no competing interests.

280

281 **Figure Legends/Captions.**

282

283 Figure 1: Fast Scoring with Indexing

284

285 a: **Similarity score.** (i) In exact search, MASST searches a query spectrum against all database spectra with similar
286 precursor masses, and computes the ExactScore, a sum of multiplications between intensities of peaks shared by the query
287 and database spectrum (shown in solid grey). In this case the score is $6.2 * 3.2 + 10.2 * 16.3 = 186.1$. (ii) In the case of
288 analog search, MASST searches the query spectrum against all database spectra within a specific precursor mass range
289 (e.g. 300 Da) and computes the ShiftedScore, a sum of multiplications between intensities of peaks that are shared and Δ -
290 shifted between query and database spectrum. Here there is one shared (solid grey) and two Δ -shifted (dashed grey) peaks,
291 yielding a total score of $6.2 * 2.2 + 10.2 * 9.2 + 15.4 * 9.2 = 249.16$. Δ denotes the precursor mass difference between
292 query and database spectra. b: **Fast Dot Product.** (i) Given a database of spectra the fast dot procedure starts with (ii)
293 constructing an indexing table, where each row corresponds to a fragment peak mass, and contains a list of tuples of
294 spectra indices that contain the peak, along with the intensity of the peak in these spectra. (iii) Given a query spectrum, all
295 lists corresponding to peaks present in the query are retrieved. Then, (iv) for each list, and for each tuple in the list, the
296 product of the intensity of the corresponding query peak and database peak is added to the total dot product score of query
297 and database spectra. For simplicity, in this illustration all the spectra have the same precursor mass. c: **Fast Dot Product**
298 **Indexing.** The fast dot product indexing table corresponds to a two-dimensional grid, with precursor mass on the x-axis
299 and peak mass on the y-axis. Each database peak is inserted into a list corresponding to a specific location in the grid,
300 determined by the peak mass and the precursor mass. In exact search, for each query peak only the list in a single cell will
301 be retrieved (highlighted with green circle). For analog search, red cells (corresponding the shared peaks) and blue cells
302 (corresponding to Δ -shifted peaks) are retrieved.

303

304 Figure 2: MASST+, Clustering+, Networking+ enables Lanthepeptide Discovery

305

306 a: **MASST+ Performance.** (i) MASST+ is two orders of magnitudes faster than MASST in exact and analog search for
307 various database sizes. (ii) MASST+ outperforms MASST in memory efficiency. b: **Clustering+ and Networking+**
308 **Performance.** (i) Clustering+ runtime versus MS-Clustering. (ii) Pairing+ runtime versus spectral networking. (iii)
309 Networking+ runtime versus Molecular Networking. Clustering+, Pairing+ and Networking+ are two order of
310 magnitudes faster than the state-of-the-art methods when processing large datasets. c: **Lanthepeptides.** (i) Biosynthetic
311 gene cluster of CHM-1731. Genes with different functions are highlighted with different colors. (ii) Annotation of peaks in
312 mass spectrum representing CHM-1731. B-ions (prefix fragmentations) are shown in blue, and y-ions (suffix
313 fragmentations) are shown in red. (iii) Mass error of annotations are shown in parts per million (ppm). Stars stand for
314 dehydrated serine / threonine.

315 **References**

316 1.da Silva, R. R., Dorrestein, P. C. & Quinn, R. A. Illuminating the dark matter in metabolomics. *Proceedings of*
317 *the National Academy of Sciences* **112**, 12549–12550 (2015).

318

- 319 2.Aron, A. T. *et al.* Reproducible molecular networking of untargeted mass spectrometry data using GNPS.
320 *Nature protocols* **15**, 1954–1991 (2020).
- 321
- 322 3.Courraud, J., Ernst, M., Svane Laursen, S., Hougaard, D. M. & Cohen, A. S. Studying autism using untargeted
323 metabolomics in newborn screening samples. *Journal of Molecular Neuroscience* **71**, 1378–1393 (2021).
- 324
- 325 4.Depke, T., Thöming, J. G., Kordes, A., Häussler, S. & Brönstrup, M. Untargeted LC-MS metabolomics
326 differentiates between virulent and avirulent clinical strains of *Pseudomonas aeruginosa*. *Biomolecules* **10**, 1041
327 (2020).
- 328
- 329 5.Eberhard, F. E., Klimpel, S., Guarneri, A. A. & Tobias, N. J. Metabolites as predictive biomarkers for
330 Trypanosoma cruzi exposure in triatomine bugs. *Computational and structural biotechnology journal* **19**, 3051–
331 3057 (2021).
- 332
- 333 6.Ernst, M. *et al.* Gestational age-dependent development of the neonatal metabolome. *Pediatric Research* **89**,
334 1396–1404 (2021).
- 335
- 336 7.Frank, A. M. *et al.* Clustering millions of tandem mass spectra. *Journal of proteome research* **7**, 113–122
337 (2008).
- 338
- 339 8.Jarmusch, A. K. *et al.* ReDU: a framework to find and reanalyze public mass spectrometry data. *Nature
340 methods* **17**, 901–904 (2020).
- 341
- 342 9.Kale, N. S. *et al.* MetaboLights: an analog-access database repository for metabolomics data. *Current
343 protocols in bioinformatics* **53**, 14–13 (2016).
- 344

- 345 10.Kuo, T.-H., Yang, C.-T., Chang, H.-Y., Hsueh, Y.-P. & Hsu, C.-C. Nematode-trapping fungi produce diverse
346 metabolites during predator-prey interaction. *Metabolites* **10**, 117 (2020).
- 347
- 348 11.Lybbert, A. C., Williams, J. L., Raghuvanshi, R., Jones, A. D. & Quinn, R. A. Mining public mass
349 spectrometry data to characterize the diversity and ubiquity of *P. aeruginosa* specialized metabolites.
350 *Metabolites* **10**, 445 (2020).
- 351
- 352 12.Mohimani, H. *et al.* Dereplication of peptidic natural products through database search of mass spectra.
353 *Nature chemical biology* **13**, 30–37 (2017).
- 354
- 355 13.Nothias, L.-F. *et al.* Feature-based molecular networking in the GNPS analysis environment. *Nature methods*
356 **17**, 905–908 (2020).
- 357
- 358 14.Petas, D. *et al.* Non-Targeted Metabolomics Enables the Prioritization and Tracking of Anthropogenic
359 Pollutants in Coastal Seawater. (2020).
- 360
- 361 15.Quinn, R. A. *et al.* Global chemical effects of the microbiome include new bile-acid conjugations. *Nature*
362 **579**, 123–129 (2020).
- 363
- 364 16.Sud, M. *et al.* Metabolomics Workbench: An international repository for metabolomics data and metadata,
365 metabolite standards, protocols, tutorials and training, and analysis tools. *Nucleic acids research* **44**, D463–
366 D470 (2016).
- 367
- 368 17.van Der Hooft, J. J. *et al.* Linking genomics and metabolomics to chart specialized metabolic diversity.
369 *Chemical Society Reviews* **49**, 3297–3314 (2020).
- 370

- 371 18.Wang, M. *et al.* Sharing and community curation of mass spectrometry data with Global Natural Products
372 Social Molecular Networking. *Nature biotechnology* **34**, 828–837 (2016).
- 373
- 374 19.Wang, M. *et al.* Mass spectrometry searches using MASST. *Nature biotechnology* **38**, 23–26 (2020).
- 375
- 376 20.Ramos, A. E. F., Evanno, L., Poupon, E., Champy, P. & Beniddir, M. A. Natural products targeting strategies
377 involving molecular networking: different manners, one goal. *Natural product reports* **36**, 960–980 (2019).
- 378
- 379 21.Kalinski, J.-C. J. *et al.* Molecular networking reveals two distinct chemotypes in pyrroloiminoquinone-
380 producing Tsitsikamma favus sponges. *Marine drugs* **17**, 60 (2019).
- 381
- 382 22.Raheem, D. J., Tawfike, A. F., Abdelmohsen, U. R., Edrada-Ebel, R. & Fitzsimmons-Thoss, V. Application
383 of metabolomics and molecular networking in investigating the chemical profile and antitrypanosomal activity
384 of British bluebells (*Hyacinthoides non-scripta*). *Scientific reports* **9**, 1–13 (2019).
- 385
- 386 23.Trautman, E. P., Healy, A. R., Shine, E. E., Herzon, S. B. & Crawford, J. M. Domain-targeted metabolomics
387 delineates the heterocycle assembly steps of colibactin biosynthesis. *Journal of the American Chemical Society*
388 **139**, 4195–4201 (2017).
- 389
- 390 24.Vizcaino, M. I., Engel, P., Trautman, E. & Crawford, J. M. Comparative metabolomics and structural
391 characterizations illuminate colibactin pathway-dependent small molecules. *Journal of the American Chemical
392 Society* **136**, 9244–9247 (2014).
- 393
- 394 25.Nguyen, D. D. *et al.* Indexing the *Pseudomonas* specialized metabolome enabled the discovery of poaeamide
395 B and the bananamides. *Nature microbiology* **2**, 1–10 (2016).
- 396

- 397 26.Woo, S., Kang, K. B., Kim, J. & Sung, S. H. Molecular networking reveals the chemical diversity of
398 selaginellin derivatives, natural phosphodiesterase-4 inhibitors from *Selaginella tamariscina*. *Journal of natural*
399 *products* **82**, 1820–1830 (2019).
- 400
- 401 27.Reginaldo, F. P. S. *et al.* Molecular Networking Discloses the Chemical Diversity of Flavonoids and
402 Selaginellins in *Selaginella convoluta*. *Planta Medica* **87**, 113–123 (2021).
- 403
- 404 28.Frank, A. M. *et al.* Spectral archives: extending spectral libraries to analyze both identified and unidentified
405 spectra. *Nature methods* **8**, 587–591 (2011).
- 406
- 407 29.Bandeira, N., Tsur, D., Frank, A. & Pevzner, P. A. Protein identification by spectral networks analysis.
408 *Proceedings of the National Academy of Sciences* **104**, 6140–6145 (2007).
- 409
- 410 30.Schnell, N. *et al.* Prepeptide sequence of epidermin, a ribosomally synthesized antibiotic with four sulphide-
411 rings. *Nature* **333**, 276–278 (1988).
- 412
- 413 31.Mohr, K. I. *et al.* Pinensins: the first antifungal lantibiotics. *Angewandte Chemie International Edition* **54**,
414 11254–11258 (2015).
- 415
- 416 32.Férir, G. *et al.* The lantibiotic peptide labyrinthopeptin A1 demonstrates broad anti-HIV and anti-HSV
417 activity with potential for microbicidal applications. *PloS one* **8**, e64010 (2013).
- 418
- 419 33.Iorio, M. *et al.* A glycosylated, labionin-containing lanthipeptide with marked antinociceptive activity. *ACS*
420 *chemical biology* **9**, 398–404 (2014).
- 421
- 422 34.Arnison, P. G. *et al.* Ribosomally synthesized and post-translationally modified peptide natural products:

- 423 overview and recommendations for a universal nomenclature. *Natural product reports* **30**, 108–160 (2013).
- 424
- 425 35.Frank, A. & Pevzner, P. PepNovo: de novo peptide sequencing via probabilistic network modeling.
- 426 *Analytical chemistry* **77**, 964–973 (2005).
- 427
- 428 36.Mohimani, H. *et al.* Dereplication of peptidic natural products through database search of mass spectra.
- 429 *Nature chemical biology* **13**, 30–37 (2017).
- 430
- 431 37.Yang, J. Y. *et al.* Molecular networking as a dereplication strategy. *Journal of natural products* **76**, 1686–
- 432 1699 (2013).
- 433
- 434 38.Ramos, A. E. F., Evanno, L., Poupon, E., Champy, P. & Beniddir, M. A. Natural products targeting strategies
- 435 involving molecular networking: different manners, one goal. *Natural product reports* **36**, 960–980 (2019).
- 436
- 437 39.Watrous, J. *et al.* Mass spectral molecular networking of living microbial colonies. *Proceedings of the*
- 438 *National Academy of Sciences* **109**, E1743–E1752 (2012).
- 439
- 440 40.Walker, M. C. *et al.* Precursor peptide-targeted mining of more than one hundred thousand genomes expands
- 441 the lanthipeptide natural product family. *BMC genomics* **21**, 1–17 (2020).
- 442
- 443 41.Kodani, S., Lodato, M. A., Durrant, M. C., Picart, F. & Willey, J. M. SapT, a lanthionine-containing peptide
- 444 involved in aerial hyphae formation in the streptomycetes. *Molecular microbiology* **58**, 1368–1380 (2005).
- 445
- 446 42.Ueda, K. *et al.* AmfS, an extracellular peptidic morphogen in *Streptomyces griseus*. *Journal of Bacteriology*
- 447 **184**, 1488–1492 (2002).
- 448

449 43. Bittremieux et al. Analog Access Repository-Scale Propagated Nearest Neighbor Suspect Spectral Library
450 for Untargeted Metabolomics. *BioRxiv*, [Preprint] (2022) Available from:
451 <https://doi.org/10.1101/2022.05.15.490691>

452

453 44. Ludwig, M., Fleischauer, M., Dührkop, K., Hoffmann, M. A. & Böcker, S. De novo molecular formula
454 annotation and structure elucidation using SIRIUS 4. *Computational Methods and Data Analysis for*
455 *Metabolomics* 185–207 (2020).

456

457 45. Dührkop, K. et al. Systematic classification of unknown metabolites using high-resolution fragmentation
458 mass spectra. *Nature Biotechnology* **39**, 462–471 (2021).

459

460 46. Mohimani, H., Kim, S. and Pevzner, P.A., A new approach to evaluating statistical significance of spectral
461 identifications. *Journal of proteome research*, **12**(4), pp.1560-1568.

462

463

464

465

466

Methods

467 **Overview of MASST algorithm.** In exact search mode, MASST performs exact search by
468 retrieving the spectra in the database that have the same precursor mass as the query and computing
469 SharedScore between each retrieved spectrum and the query. Analog search is conducted by
470 retrieving all spectra within a large precursor mass tolerance (e.g. 300 Da) of the query precursor
471 mass, and computing the ShiftedScore (Figure 1a-ii). To compute these scores, MASST iterates
472 over all the peaks in the query spectrum, and for each peak it explores whether a peak with similar
473 or shifted m/z is present in each database spectrum. Whenever such a peak is present, MASST

474 increments the score between the query and that database spectrum by the product of the intensity
475 of peaks in the query and the database spectrum.

476 **MASST+ exact search.** Given a query spectrum, MASST+ efficiently searches a database of
477 reference spectra to find similar spectra by using the fast dot product algorithm (Figure 1b). For
478 each precursor mass M and each peak mass p , a list of indices of all spectra with precursor M and
479 peak within a tolerance threshold of p are stored, along with intensity of peaks. In case of exact
480 search, given a query spectrum with precursor mass M , MASST+ iterates through the peaks in the
481 query spectrum and retrieves the lists corresponding to the peaks and precursor mass M . As each list
482 is stored on disk, each list can be retrieved in $O(1)$ time. The SharedScore is then calculated by
483 multiplying and adding up the intensity of each peak in the query spectrum and reference spectra
484 (Figure 1b-iv).

485 **MASST+ analog search.** In the case of analog search, MASST+ uses a large precursor mass
486 tolerance (e. g. 300Da) and computes ShiftedScore (Figure 1a-ii). ShiftedScore takes into account
487 both shared and Δ -shifted peaks, where Δ is the mass difference between the query and each
488 reference spectrum. In analog mode, all reference spectra are processed into lists as in MASST+
489 exact search. Given a query spectrum, MASST+ analog search iterates through each peak \square in the
490 query spectrum with precursor mass \square , and scans lists (\square', \square') where either $\square = \square'$ (shared peak) or
491 $\square - \square = \square' - \square'$ (shifted peak). The ShiftedScore between the query and each reference spectrum is
492 calculated by multiplying and adding up the intensity of shared and shifted peaks in the two spectra
493 (Supplementary Fig. 4). Note that MASST+ analog search is a variant of the fast dot product
494 algorithm (Figure 1b) as both methods rely on similarly structured index tables. Rather than just
495 retrieving one list for each query spectrum peak, however, MASST+ analog search retrieves two
496 lists.

497 **MASST+ indexing.** To handle continuous values of peak masses, we bin peak masses into discrete
498 values. Depending on the bin size and product mass tolerance, one or more bins must be retrieved

499 when processing each query peak during search. We use a bin size of 0.01Da, which can handle
500 both high-resolution (0.01Da accuracy) and low-resolution (0.5Da accuracy) data.

501

502 **Overview of Molecular Networking.** In order to find structurally related families of small
503 molecules, the existing molecular networking method first clusters spectra from identical molecules
504 using MS-Clustering²⁸. It then connects clusters of related molecules using spectral networking²⁹.
505 MS-Clustering puts two spectra in the same cluster if their precursor mass difference is below a
506 threshold (usually 2 Da) and their cosine dot product (a normalized SharedScore) is above a certain
507 threshold (usually 0.7). Then for each cluster, a consensus spectrum is constructed using the
508 approach introduced by Frank et al²⁸. In spectral networking, two consensus spectra are connected
509 to each other if the shared-shifted cosine score (normalized ShiftedScore) is above a threshold
510 (default is 0.7).

511

512 **Networking+ algorithm.** Networking+ consists of two modules, Clustering+ and Pairing+.
513 Clustering+ is implemented using a greedy procedure (Supplementary Fig. 12). Given a dataset of
514 N spectra, Clustering+ creates an initial cluster whose center is set to be the first spectrum in the
515 dataset. Then in the following N-1 iterations, the similarity score between each remaining spectra
516 and all the existing cluster centers is calculated. To efficiently calculate the similarity score between
517 the spectrum and all cluster centers, an indexing table similar to MASST+ exact search is
518 constructed and iteratively updated. For each precursor mass M and peak mass p , the indexing table
519 stores the list of all clusters that have centers with a specific precursor mass M and a peak mass \square .
520 At each iteration, whenever the highest score between the spectrum and cluster centers is greater
521 than a threshold (default is 0.7), the spectrum is added to the highest-scoring cluster, and the center
522 of the cluster is updated. If the highest score is below the threshold, then a new cluster is created,
523 and the current spectrum is set as the center of the cluster. This procedure continues until all the
524 spectra are clustered.

525 To maintain efficiency, whenever a new spectrum is added, the center is updated only when the
526 cluster size doubles (e.g. after the addition of the first, second, fourth, eighth, sixteenth, etc.
527 spectrum to the cluster). Similar to Frank et al²⁸, the center is computed by adding peaks that are
528 present in the majority of the members of the cluster. The intensity of each peak is calculated as the
529 average of the intensity of the corresponding peaks in members. All spectra are initially normalized.

530 Pairing+ computes a score similar to MASST+ analog search (Supplementary Fig. 4) that accounts
531 for Δ -shifted and shared peaks for all pairs of input spectra (e.g. cluster centers from clustering+).
532 To do this, it constructs an indexing table similar to MASST+ analog search. Then the table is used
533 to efficiently compute the score between all pairs of spectra (Supplementary Fig. 13).

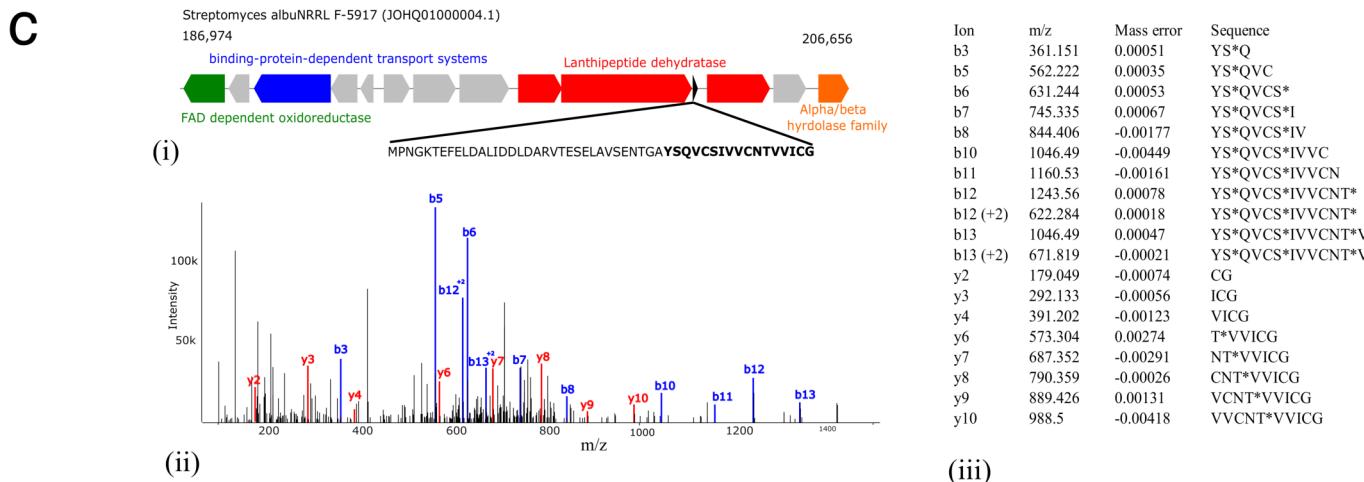
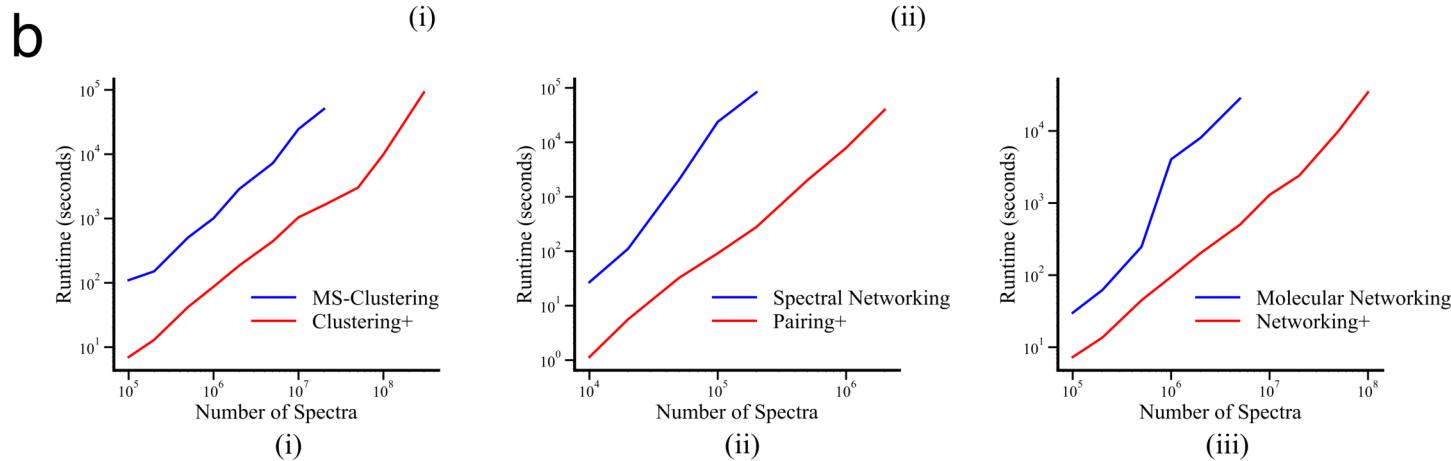
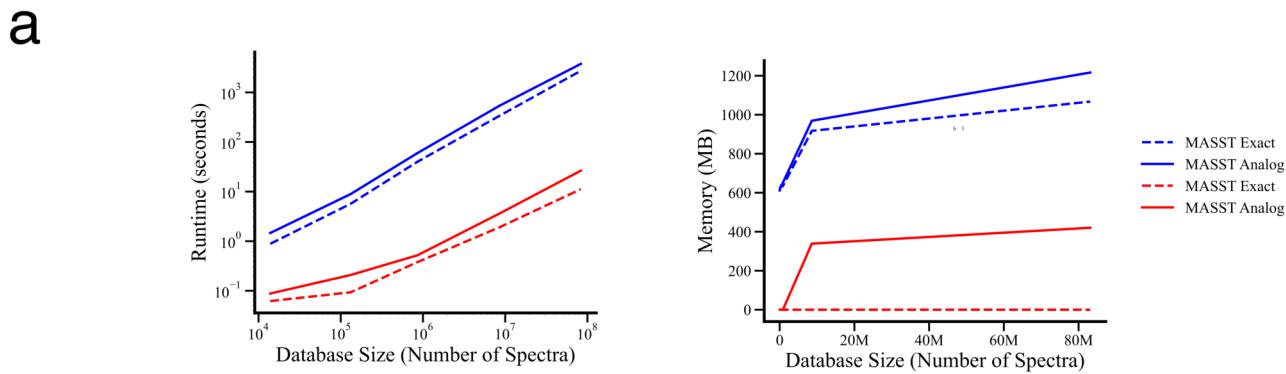
534

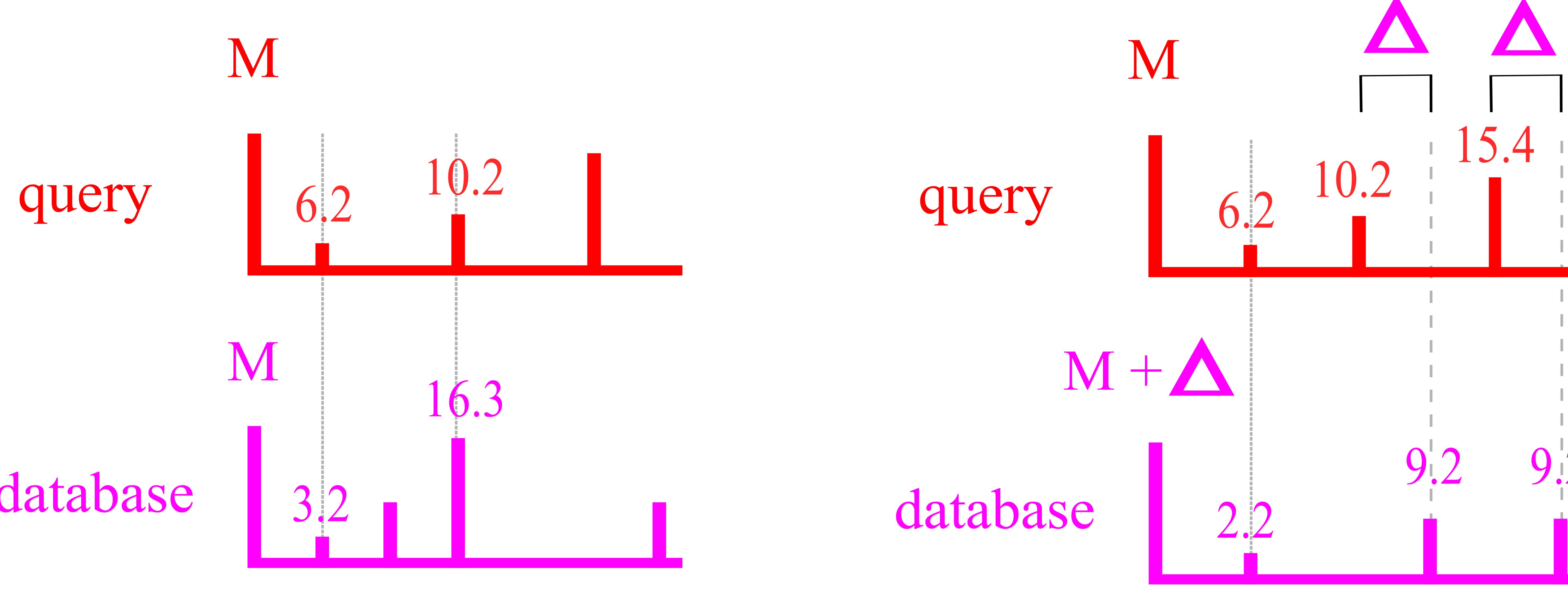
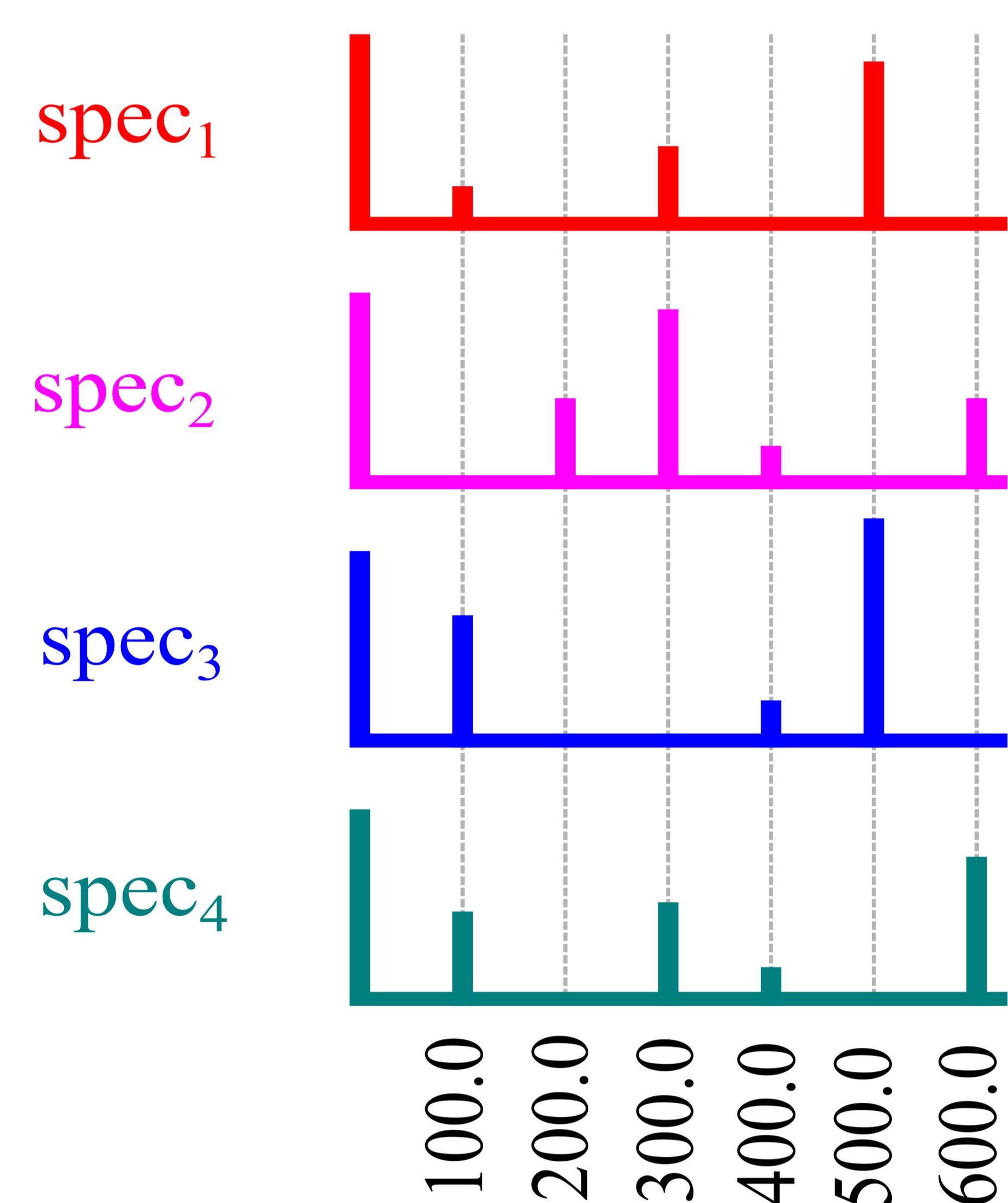
535

536 **Data Availability.** The datasets analyzed are available at gnps.ucsd.edu. Accession codes related
537 to Lanthepeptide portion of manuscript are MSV000090476 ,MSV000090473 ,MSV000090472
538 ,MSV000090471 ,MSV000090457 ,MSV000089818 ,MSV000089817 ,MSV000089816
539 ,MSV000089815 ,MSV000089813 ,MSV000088816 ,MSV000088801 ,MSV000088800
540 ,MSV000088764 ,MSV000088763 . For comparing MASST+ and Networking+ against previous
541 state of the art, datasets MSV000078787 , Clustered GNPS, and Unclustered GNPS were used. The
542 Accession codes for Clustered GNPS and Unclustered GNPS are in Supplementary_Data_1.xlsx.

543

544 **Code Availability.** MASST+ and Networking+ presented in the paper are available at
545 <https://github.com/mohimanilab/MASSTplus>. Other custom software utilized in the paper include
546 Seq2Ripp (<https://github.com/mohimanilab/seq2ripp>), PepNovo
547 (<https://github.com/jmchilton/pepnovo>), and Dereplicator ([https://ccms-
548 \[ucsd.github.io/GNPSDocumentation/dereplicator/\]\(https://ccms-ucsd.github.io/GNPSDocumentation/dereplicator/\)](https://ccms-ucsd.github.io/GNPSDocumentation/dereplicator/)).

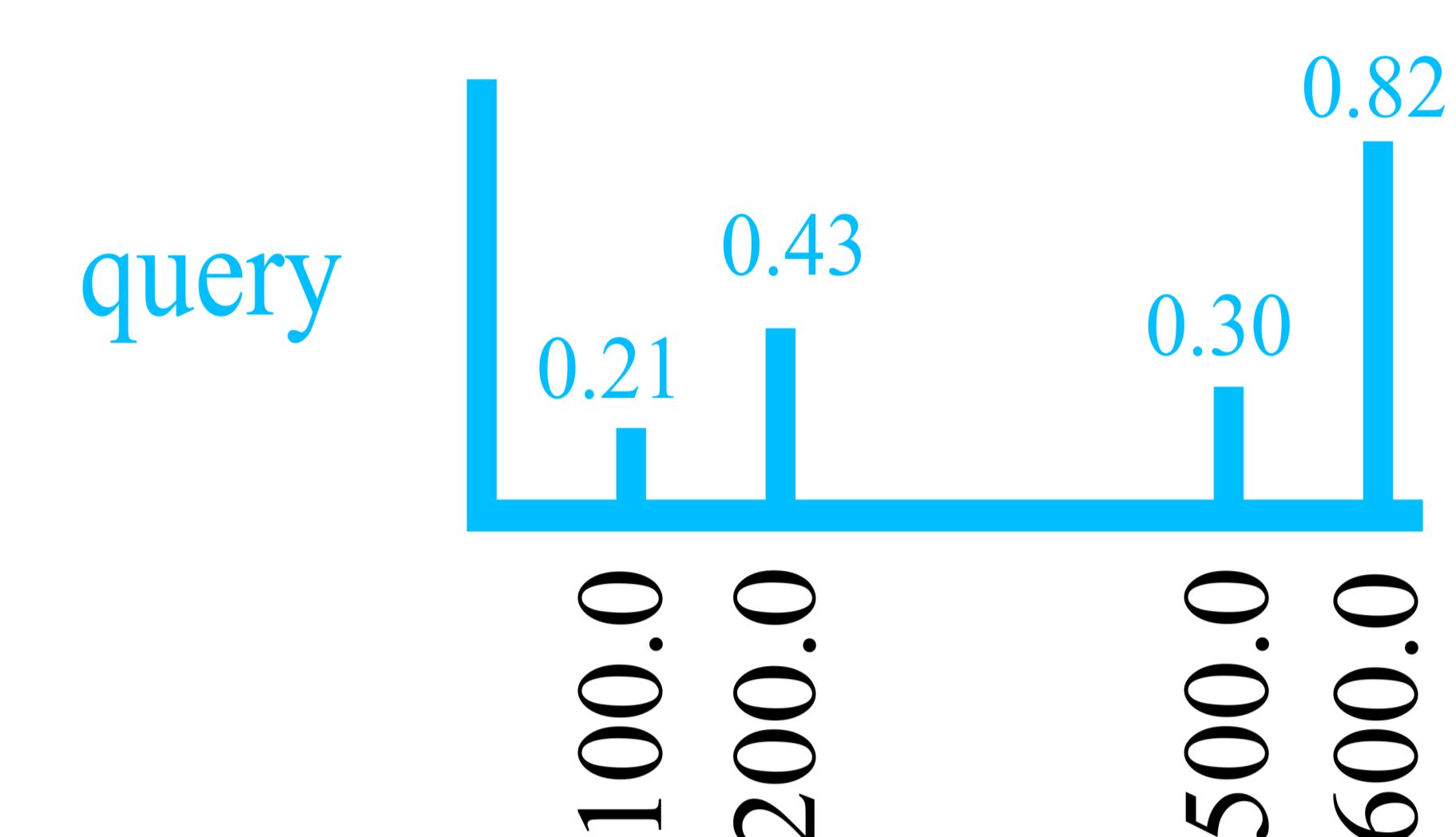


a**b**

(i) reference spectra

Peak	(SpecIndex, Intensity)
100.0	{(1,0.23),(3,0.51),(4,0.45)}
200.0	{(2,0.40)}
300.0	{(1,0.43),(2,0.80),(4,0.50)}
400.0	{(2,0.19),(3,0.18),(4,0.19)}
500.0	{(1,0.87),(3,0.84)}
600.0	{(2,0.40),(4,0.72)}

(ii) constructing the indexing table



100.0	{(1,0.23),(3,0.51),(4,0.45)}
200.0	{(2,0.40)}
500.0	{(1,0.87),(3,0.84)}
600.0	{(2,0.40),(4,0.72)}

(iii) retrieving the lists corresponding to peaks in the query spectrum

	process 100.0	process 200.0	process 500.0	process 600.0	score
Score(query,spec1) =	$0.23 * 0.21$	+	$0.87 * 0.30$	+	= 0.31
Score(query,spec2) =		+	$0.40 * 0.43$	+	= 0.50
Score(query,spec3) =	$0.51 * 0.21$	+	$0.84 * 0.30$	+	= 0.36
Score(query,spec4) =	$0.45 * 0.21$	+		+	$0.72 * 0.82$ = 0.69

(iv) calculating the scores

c