

1 **Title:** Annotating full-scan MS data using tandem MS libraries

2

3

4 **Authors:**

5 Shipei Xing^{1,2}, Vincent Charron-Lamoureux^{1,2}, Yasin El Abiead^{1,2}, Pieter C Dorrestein^{1,2,3,4,*}

6

7

8 **Author affiliations:**

9 ¹Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California San Diego,
10 La Jolla, CA, 92093, USA

11 ²Collaborative Mass Spectrometry Innovation Center, University of California San Diego, La Jolla,
12 CA, 92093, USA

13 ³Department of Pharmacology, University of California San Diego, La Jolla, CA, 92093, USA

14 ⁴Center for Microbiome Innovation, University of California San Diego, La Jolla, CA, 92093, USA

15

16 *Author to whom correspondence should be addressed.

17

18

19

20

21

22

23 **Abstract (70 words, currently 70)**

24

25 Full-scan mass spectrometry (MS) data from both liquid chromatography (LC) and MS
26 imaging capture multiple ion forms, including their in-source fragments. Here we leverage such
27 fragments to structurally annotate full-scan data from LC-MS or MS imaging by matching against
28 peak intensity scaled tandem MS spectral libraries using precursor-tolerant reverse match
29 scoring. Applied to inflammatory bowel disease and imaging datasets, we show the approach
30 facilitates re-analyses of data in public repositories.

31

32 **Main (1500 words, currently 2807)**

33

34 The scientific community is investing substantial resources – funds, time, and effort – in
35 collecting samples, obtaining metabolomics data, analyzing results, and making them publicly
36 available. Just as a conservative number, we are approaching a milestone of one million public
37 liquid chromatography-mass spectrometry (LC-MS) files and over 11,700 MS imaging files across
38 major repositories such as Metabolomics Workbench¹, MetaboLights², GNPS/MassIVE³, and
39 METASPACE⁴. The rationale behind public data deposition (even when restricted access is
40 employed) extends beyond promoting scientific transparency and reproducibility; it also facilitates
41 future reuse, as typically only a small fraction of the data is utilized upon initial publication.

42

43 Given that the average annotation rate in metabolomics studies ranges from 10% to 20%⁵⁻
44 ⁸, one crucial aspect of data reuse is to provide new annotations that can be re-contextualized to
45 uncover new biological insights. In untargeted metabolomics, two primary approaches are
46 employed: data collection with tandem MS (MS/MS) and without (MS1-only or full-scan). The
47 latter method offers an advantage in peak shape quality due to the increased number of scans
48 contributing to each peak, leading to enhanced absolute or relative quantification accuracy and
49 thus more reliable statistical analyses. However, MS1-only data presents limitations in discovering
50 molecules that were detected but not yet annotated through subsequent reanalysis. Notably, more
51 than 40% of untargeted LC-MS metabolomics data files in public repositories consist solely of
52 MS1 data, and almost all MS imaging metabolomics data are MS1 scans. This situation creates
53 a significant gap in data reuse and reinvestigation for full-scan data in untargeted metabolomics.

54

55 Traditionally, MS1 data interpretation relies on accurate mass measurements and isotopic
56 patterns, which can suggest possible molecular formulas^{4,9,10} but often falls short of providing
57 structural information. However, it is also generally understood that many ions may undergo in-
58 source fragmentation or exhibit post-source decay¹¹⁻¹⁶, generating fragment pieces that also
59 appear in MS1 data. As these processes involve thermal activation, the resulting in-source
60 fragment ions exhibit fragmentation patterns very similar to those observed in collision-induced
61 dissociation (CID) MS/MS spectra. This opens the possibility of leveraging such in/post-source
62 fragments as a handle to create pseudo MS/MS spectra, also referred to as composite spectra¹⁷,
63 that can be leveraged for MS/MS reference library-based annotation in metabolomics and
64 exposomics studies¹⁵⁻²³.

65

66 Strategies, such as IDSL.CSA¹⁷, have demonstrated the proof-of-principle of matching
67 pseudo MS/MS spectra, obtained by aggregating ion forms across entire datasets, to reference
68 MS/MS libraries using scoring methods like cosine or entropy similarity²⁴. This works particularly
69 well for GC-MS²⁵ due to the consistent use of a fixed energy (70 eV) for both data acquisition and
70 reference spectra, and the absence of many co-eluting ion forms such as different adducts or

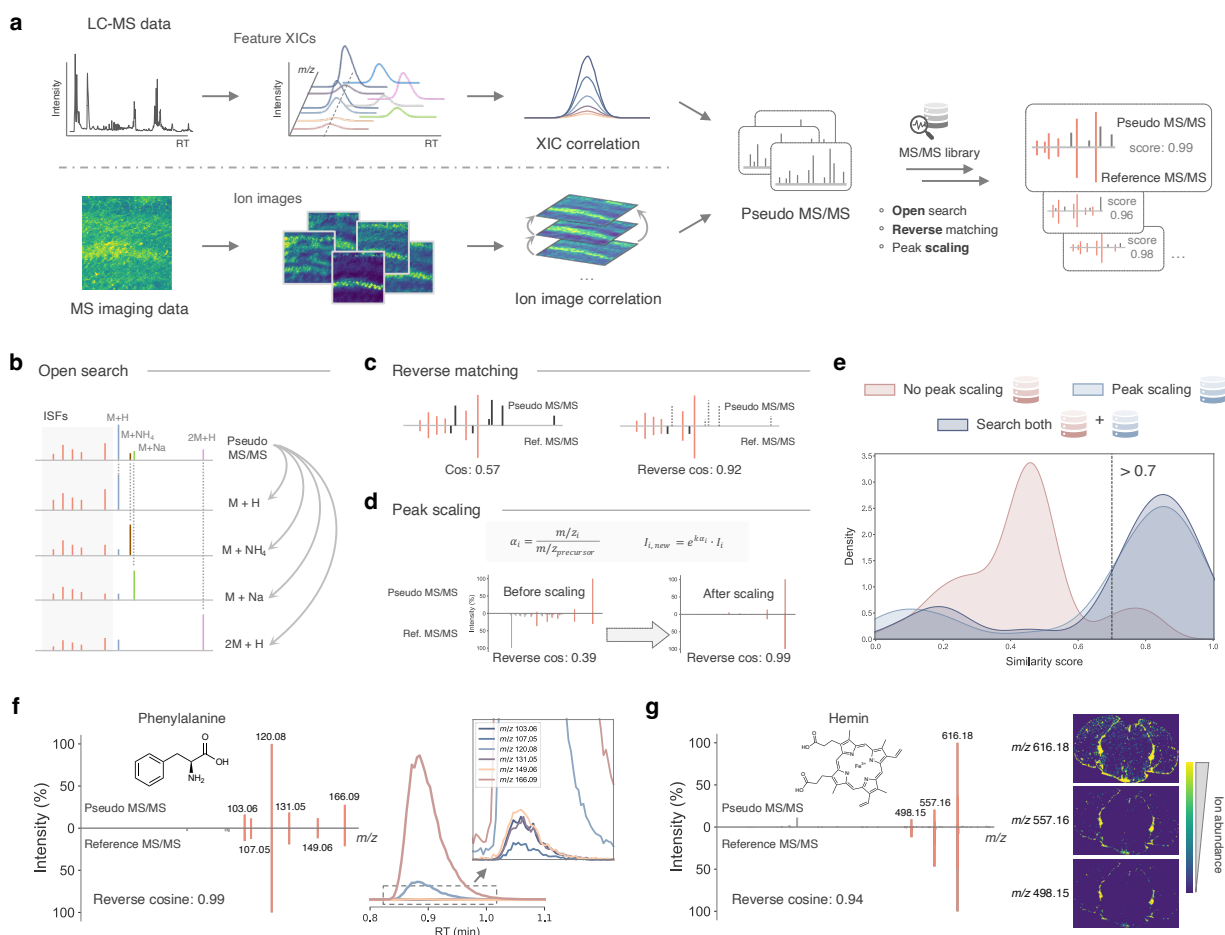
71 multi-meric species. However, in LC-MS, different ion forms, such as adducts and multimers,
72 often dominate pseudo MS/MS spectra, which may prevent matching to reference MS/MS
73 libraries that do not account for these ion forms. Furthermore, as we show below, the fragment
74 ions we detected tended to match reference MS/MS spectra generated with lower-energy
75 fragmentation, while most reference spectra available in the public domain are collected under
76 medium to high collision energies. In these lower-energy spectra, low-*m/z* ions often appear at
77 very low intensities or may be absent entirely, resulting in missed matches. Other experimental
78 strategies aimed to overcome this limitation—for example, eISA²⁶ and EISA-EXPOSOME¹⁸ have
79 been developed to incorporate in-source fragments in metabolite annotation. Therefore, in their
80 current implementations, these methods require full-scan data under enhanced in-source CID
81 (isCID) energies (e.g., 40 eV) to obtain adequate spectral matches, and they work in a targeted
82 fashion. Consequently, these methodologies cannot be used for reanalyzing hundreds of
83 thousands of public untargeted MS1 data files acquired without extra isCID experimental designs.

84
85 In the realm of MS imaging, metabolite candidate annotations are obtained by annotating
86 molecular formulas and cross-referencing these formulas in common metabolite/lipid chemical
87 databases such as HMDB²⁷, LipidMaps²⁸ and ChEBI²⁹. This approach corresponds to level 4
88 identification confidence (unequivocal molecular formula) according to the Metabolomics
89 Standards Initiative³⁰. Despite the rapid growth of MS/MS spectral libraries⁸, metabolite annotation
90 of MS imaging data has not yet fully benefited from this expanding community resource.

91
92 We therefore highlight in this work that we can annotate MS1 data applicable to both LC-
93 MS and MS imaging data. Our approach integrates two steps (although how these steps are
94 implemented is critical): (1) clustering ions or metabolic features through correlation analysis of
95 extracted ion chromatograms (XICs) or ion images in the retention time domain or spatial manner
96 (**Fig. 1a**); and (2) employing a precursor-tolerant (open search³¹) but using reverse spectral
97 matching approach to compare deconvolved MS1 spectra, or pseudo MS/MS spectra, against
98 peak scale-adjusted reference MS/MS libraries for structure candidate identification. Unlike
99 traditional forward spectral matching which utilizes all the peaks in both query and reference
100 spectra for scoring, reverse matching is a unidirectional spectral comparison which ignores
101 unaligned peaks in the query spectrum³², tolerating contaminant peaks sourced from co-eluting
102 metabolic features or signal artifacts. In the following sections, we elaborate on the spectral
103 matching design and its underlying rationales, demonstrating how this approach enhances
104 annotation capabilities for MS1 data in both LC-MS and MS imaging experiments.

105
106 Each molecule detected by mass spectrometry is represented by multiple molecular ions
107 of various adduct forms (e.g., $M+NH_4^+$, $M+Na^+$) that co-elute during chromatography^{33–35} or, in the
108 case of MS imaging, share spatial correlations with each other, in addition to in-source fragments.
109 Consequently, intact molecular ions of different adducts, along with their fragments, appear in the

110 same reconstructed pseudo MS/MS spectrum (**Fig. 1b**). Unlike a typical MS/MS spectrum
 111 collected in data-dependent acquisition (DDA) mode, it is not known which ion—if any—represents
 112 the precursor ion in a pseudo MS/MS. To address this, we implemented an open search
 113 approach—it does not assume a single, predefined precursor ion for each spectrum, but instead
 114 considers every ion as a potential precursor ion simultaneously. It employs an unlimited mass
 115 tolerance window to accommodate potential mass shifts due to different adducts or multimers,
 116 enabling the recognition of various precursor types within the same spectrum.
 117



118 **Fig. 1** | Structure annotation of full-scan MS data. **a**, A unified solution to annotate MS1 data from LC-MS or MS imaging experiments.
 119 Pseudo MS/MS spectra are generated through correlation analyses in time or spatial domain. The precursor ion-tolerant (open search)
 120 reverse spectral searching allows structure annotation of pseudo MS/MS leveraging existing reference MS/MS libraries. **b**, Open
 121 search allows matching against reference spectra of multiple adduct forms. ISFs, in-source fragments. **c**, Reverse spectral search
 122 discards unmatched peaks in the query spectrum and thus improves spectral search. **d**, Peak intensity scaling helps to match pseudo
 123 MS/MS spectra against reference MS/MS which are collected under medium to high collision energies. **e**, Similarity score distribution
 124 of searching pseudo MS/MS against libraries with or without peak scaling, or both. Reverse cosine scores of ground truths using
 125 chemical standards are used. **f**, An example of structure annotation from LC-MS data (NIST human feces) and XIC correlations. Peak
 126 intensities are square rooted. **g**, An example of structure annotation from MS imaging data (mouse brain) and extracted ion images.
 127 Ion images were created using 5 ppm mass tolerance. Peak intensities are square rooted.
 128
 129

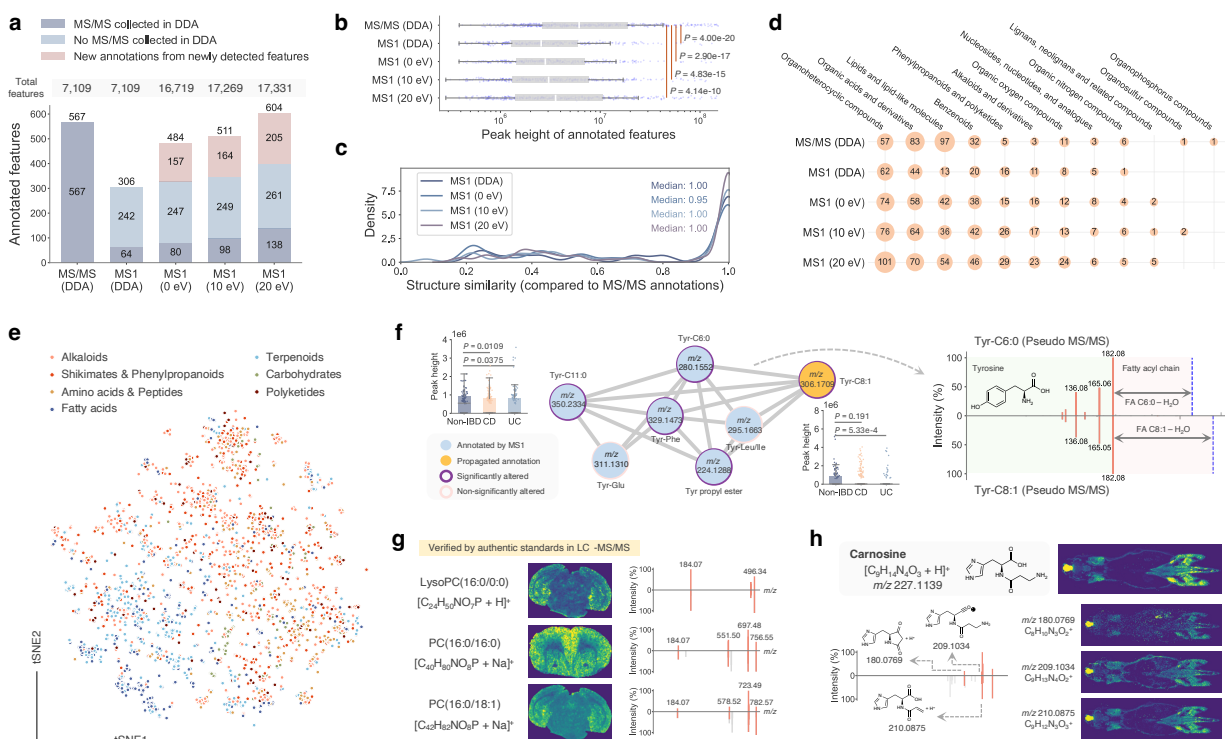
130 However, there are challenges in matching a pseudo MS/MS spectrum against the MS/MS
131 spectral library, making it often not possible to provide direct matches. Pseudo MS/MS spectra
132 contain not only fragment ions but also molecular ions of other adduct types and unavoidably mis-
133 clustered ions that are co-eluting. These additional ions are undesirable during spectral matching
134 as they significantly diminish the search scores as the reference libraries do not contain all of the
135 different ion forms. We therefore employed a reverse spectral search (**Fig. 1c**). In this method,
136 reference spectra serve as templates, and unmatched peaks in the pseudo MS/MS are discarded
137 when calculating matching scores. This approach is particularly crucial for MS imaging data,
138 where ions with similar biological functions tend to have similar spatial patterns and thus high
139 correlations, resulting in more ions that should not be compared when trying to annotate (e.g.,
140 lipid molecules can exhibit similar spatial distributions on cellular membranes).

141
142 Furthermore, as pseudo MS/MS spectra are obtained with minimal energy input (only
143 energy for transfer and/or trapping of the ions), the fragment ion intensities tend to more closely
144 align with low-energy CID spectra. Currently, most reference MS/MS spectra in libraries are
145 collected under medium to high collision energies. Therefore, we developed a peak intensity
146 scaling approach to better align them (**Fig. 1d**). Using chemical standard pools of bile acids and
147 drugs, in total containing 14 known molecules, for which full-scan MS data were collected under
148 in-source CID (isCID) energies of 0 eV, 10 eV and 20 eV, we demonstrated that this peak scaling
149 approach provided more matching scores of >0.7 for ground truths compared to not applying peak
150 scaling (**Fig. 1e**). Combining search results from both original and peak-scaled reference libraries
151 yielded the highest number of matches with reverse cosine scores larger than 0.7.

152
153 To further validate our approach, we collected LC-MS data from NIST reference human
154 fecal samples in both DDA and full-scan modes. Full-scan data were acquired under 0 eV, 10 eV,
155 and 20 eV isCID energies. We were able to obtain spectral library matches for 567, 306, 484, 511
156 and 604 metabolic features in MS/MS (DDA, 42 eV), MS1 (DDA), MS1 (0 eV), MS1 (10 eV), and
157 MS1 (20 eV) modes, respectively (**Fig. 2a**). Unexpectedly, MS1 annotation revealed a unique
158 chemical space, with the majority of annotated features in MS1 data being distinct from MS/MS
159 annotations. More than 79% of the features annotated via pseudo MS/MS lacked corresponding
160 MS/MS spectra in DDA experiments. While DDA typically acquires MS/MS spectra for the more
161 abundant features, this approach captures more low-intensity features when they produce
162 sufficient in-source fragments (**Fig. 2b**). When examining the same metabolic features collected
163 in DDA, structure similarity analyses between MS1 annotations and MS/MS annotations showed
164 that they generated similar chemical candidates (**Fig. 2c**), where a higher isCID energy led to
165 more similar or identical structure matches with MS/MS annotations. We then investigated the
166 compound classes³⁶ of annotated compounds under different acquisition conditions (**Fig. 2d**).
167 While MS1 data generally annotated more molecules than MS/MS across most compound
168 classes, organic acids & derivatives, and lipids & lipid-like molecules were not as well recognized

169 in MS1 annotation compared to MS/MS, and this suggests that certain classes of compounds will
 170 be easier to annotate via the pseudo MS/MS strategy forwarded here. Overall, above results
 171 indicate that MS1 annotation expands the range of detectable metabolites, potentially uncovering
 172 previously overlooked compounds in untargeted metabolomics studies, including those available
 173 in public repositories.

174
 175



176
 177
 178
 179
 180
 181
 182
 183
 184
 185
 186
 187
 188
 189
 190
 191
 192
 193
 194

Fig. 2 | MS1 structure annotation provides new insights. **a**, Annotations of NIST fecal samples by MS/MS (DDA) or MS1 data with different isCID energies. **b**, MS1 annotation is able to capture low-abundant metabolic features compared to MS/MS annotations in DDA. Boxes cover the interquartile range (IQR), with medians labeled. The upper whisker represents $Q3 + 1.5 \times IQR$; the lower whisker represents $Q1 - 1.5 \times IQR$. P values were calculated using two-sided Mann-Whitney U tests. **c**, Structure similarity distributions between MS1 annotations and MS/MS annotations when annotating the same metabolic features. **d**, Compound class distributions of the metabolites annotated in the NIST feces dataset. **e**, t-SNE visualization of MS1 annotations in the IBD dataset. Nodes are colored by compound pathways from NPClassifier. **f**, An example molecular network generated using pseudo MS/MS spectra in the IBD dataset (HILIC positive). Tyrosine-related compounds were annotated and linked. P values were calculated using two-sided Mann-Whitney U tests. The mirror plot shows pseudo MS/MS spectra from Tyr-C6:0 and Tyr-C8:1. **g**, LysoPC(16:0/0:0), PC(16:0/16:0) and PC(16:0/18:1) annotated in the mouse brain imaging data. They were all verified by authentic standards in LC-MS/MS data. Ion images were created using 5 ppm mass tolerance. Peak intensities are square rooted in the mirror plots, with matched peaks shown in the pseudo MS/MS. **h**, Carnosine annotated in the mouse body imaging data. Ion images were created using 5 ppm mass tolerance. Peak intensities are square rooted in the mirror plot, with matched peaks shown in the pseudo MS/MS.

To highlight reanalysis of a public project with MS information, we revisited a public LC-MS full-scan dataset from an inflammatory bowel disease (IBD) study³⁷. This dataset comprised

195 546 stool samples from three diagnostic groups: non-IBD (n = 134), Crohn's disease (CD, n =
196 266), and Ulcerative colitis (UC, n = 146). The analysis employed four distinct LC-MS modes:
197 HILIC positive, HILIC negative, C8 positive, and C18 negative. We performed MS1 annotations
198 in batches, obtaining 3010, 293, 227 and 636 unique metabolites (unique InChIKey strings) in the
199 four modes, respectively. Altogether, we identified 3802 unique metabolites with level 2/3
200 confidence³⁰. A t-SNE visualization (**Fig. 2e**), color-coded by compound pathways from
201 NPClassifier³⁸, revealed distinct clustering patterns among annotated compounds. Alkaloids and
202 shikimates & phenylpropanoids form large, spread-out clusters, suggesting diverse structural
203 variations and their prominence in the gut metabolome. Fatty acids and terpenoids form relatively
204 distinct clusters, indicating unique intensity profiles for lipid-based metabolites across the IBD
205 sample set. We constructed a molecular network using the pseudo MS/MS spectra from the HILIC
206 positive mode. A subnetwork for tyrosine-related compounds was highlighted in **Fig. 2f**, including
207 N-acyl amides that showed alterations in non-IBD vs CD or non-IBD vs UC comparisons. The
208 annotation of Tyr-C8:1 could be propagated through modified cosine-based MS/MS similarity and
209 delta masses with neighbor nodes, and the mirror plot clearly shows the fragmentation pattern of
210 tyrosine as well as the neutral loss of the fatty acyl chain. These findings align with previous IBD
211 studies³⁹⁻⁴¹, which identified alterations in lipid metabolism and N-acyl amide profiles as key
212 factors in IBD pathogenesis, and highlight that our reanalysis approach can assist in uncovering
213 clinically relevant metabolic signatures the original depositors did not describe.

214
215 To demonstrate the efficacy of our strategy on MS imaging data, we applied it onto mouse
216 brain⁴, mouse body and human hepatocytes⁴² datasets. In the mouse brain sample, we annotated
217 hemin and phosphatidylcholine (PC) lipids of varying chain lengths. **Fig. 1g** displays the ion
218 images of the hemin cation and its in-source fragments, and the visual inspection clearly revealed
219 that the ion image of the hemin cation exhibits spatial patterns highly similar to those of its in-
220 source fragments, with the expected lower abundance of in-source fragments. **Fig. 2g** illustrates
221 the annotations of LysoPC(16:0/0:0), PC(16:0/16:0) and PC(16:0/18:1), which were all verified by
222 authentic standards in LC-MS/MS data⁴. In the mouse body dataset, we obtained 143 candidate
223 annotations. Notably, carnosine was found to be localized to the brain and muscle tissues (**Fig.**
224 **2h**), aligning with its dual role as a neuroprotector and a muscle performance enhancer⁴³. In the
225 brain, carnosine's presence suggests its involvement in neurotransmitter regulation and synaptic
226 plasticity, processes crucial for learning and memory⁴⁴. In muscle tissue, it functions as an
227 intracellular buffer, regulating pH levels during physical activity, and exhibits antioxidant
228 properties⁴⁵ that may aid in recovery from exercise-induced stress. The significant abundance of
229 carnosine in these tissues underscores its importance in both neurocognitive function and
230 physical performance. Extending our analysis to a single cell analysis data set of human cell lines,
231 we examined an MS imaging dataset from differentiated human hepatocytes⁴² revealed various
232 lipid classes including phosphatidylcholines, diacylglycerols, and triacylglycerols. As an illustrative
233 example, **Extended Data Fig. 1** showcases the annotation of a diacylglycerol species. This result

234 highlights the ability of this approach to annotate complex lipids that are interpreted at the single
235 cell level. These findings collectively demonstrate the versatility of our approach across different
236 types of data, from tissue-level imaging to single-cell analysis. By enabling confident annotation
237 of molecular species in various biological contexts, our method promises to enhance our
238 understanding of spatial metabolomics and lipidomics in health and disease.

239
240 Despite its capacity to annotate MS1 data from both LC-MS and MS imaging experiments,
241 there are a number of important limitations one has to consider when applying this approach. This
242 approach will not be able to distinguish most isomers, particularly in complex metabolite mixtures
243 with inadequate chromatographic separation. These structurally similar compounds often co-elute
244 and produce similar fragments, impeding the creation of clean pseudo MS/MS spectra and their
245 subsequent distinction. This issue is notably evident in lipid analysis—molecules of the same lipid
246 class share identical characteristic fragments (e.g., the head group ion of phosphatidylcholines),
247 where integration of heuristic rules for retention orders may provide deeper insights. Currently
248 scaling is optimized for maximum number of annotations but this also results in increased
249 incorrect matches compared to no scaling (**Extended Data Fig. 2**). We expect future optimization
250 of scaling can further improve the annotation confidence. Another consideration is the potential
251 for ion contamination or incorrect ion clustering when generating pseudo MS/MS spectra,
252 especially in MS imaging data lacking chromatographic separation. Such limitations elevate the
253 risk of incorrect matches to reference MS/MS spectra. As we show, certain compound classes
254 are underrepresented (e.g., organic acids & derivatives and lipids & lipid-like molecules) in MS1
255 annotation. This underrepresentation arises from insufficient generation of in-source fragments
256 due to the comparatively low energy imparted on the ions in MS1-only scans. These constraints
257 highlight avenues for future research, including the advancement of more precise MS1 data
258 deconvolution techniques, incorporation of additional orthogonal data for isomer differentiation,
259 and refinement of spectral search algorithms specifically tailored for MS1 data annotation.

260
261 Our MS1 annotation approach unveils exciting new prospects for untargeted
262 metabolomics data reuse and analysis. A key opportunity lies in developing an MS1-based
263 MASST^{46,47} (Mass Spectrometry Search Tool) to perform reverse metabolomics⁴⁸ on LC-MS and
264 MS imaging data, which allows the contextualization of molecules (known or unknown) driven by
265 metadata integration⁴⁹ including body distributions, producing organisms, health conditions and
266 interventions. While the current MASST enables searching MS/MS spectra against public data
267 repositories using forward (modified) cosine to retrieve valuable metadata for new biology
268 discovery, MASST could now potentially be extended to the MS1 level. As a proof-of-principle,
269 we queried the pseudo MS/MS spectrum of phenylalanine-C3:0 from the NIST feces sample,
270 which was more abundant in the omnivore group than the vegan group, against the pseudo
271 MS/MS spectra pool from the IBD dataset. This search returned an MS/MS match with cosine
272 score of 0.90 (**Extended Data Fig. 3**). The matched pseudo MS/MS was also annotated as

273 phenylalanine-C3:0 in the IBD dataset, showing statistical significance in both non-IBD vs CD and
274 non-IBD vs UC comparisons with higher abundance in the non-IBD group. This indicates the
275 feasibility of MS1-based MASST across all four major repositories. Our MS1 annotation
276 approach's ability to identify low-abundance features suggests the possibility of achieving broader
277 metabolome coverage through MS1-based molecular networking³. This approach could catalyze
278 the propagation of annotations through spectral similarity analysis, revealing previously
279 unidentified metabolites and facilitating the creation of pseudo MS/MS-based suspect libraries⁵⁰
280 for future data reuse and reanalysis. With over 14,800 untargeted metabolomics datasets (~one
281 million data files) currently available in public repositories, this represents an untapped resource
282 for exploring the dark metabolome⁵—including those elusive metabolites that have thus far
283 escaped identification. As we refine and extend our MS1 annotation techniques, we anticipate an
284 extensive deepening of our understanding of complex metabolic processes and their roles in
285 diverse biological systems and disease states.
286

287 **Methods**

288

289 ***Pseudo MS/MS spectra generation***

290 For LC-MS data, metabolic features are extracted using the MassCube backend⁵¹, which
291 is a Python-based framework for untargeted metabolomics. For each pair of metabolic features
292 within the same retention time window (e.g., ± 1.5 s), peak-peak correlation is calculated using
293 their chromatographic profiles. To perform the correlation analysis between two ions, they must
294 share at least 4 consecutive MS1 scans in their chromatographic profiles.

295

296 Pseudo MS/MS spectra are then generated as follows: For each metabolic feature (target
297 feature), all other features with correlations exceeding a predefined threshold (e.g., Pearson
298 correlation coefficient ≥ 0.80) are collected. These correlated features are compiled into a pseudo
299 MS/MS spectrum for the target feature. Peak heights of the correlated features in the original MS1
300 data are used as their respective intensities in the pseudo MS/MS spectrum. Peaks that are
301 determined as isotope peaks by MassCube are excluded from pseudo MS/MS generation.

302

303 For MS imaging data, the process is adapted to account for spatial information. Each MS
304 scan undergoes noise reduction using a moving average algorithm. Within a moving window of
305 100 Da, the baseline is determined as 5 times the mean intensity of the lowest 5% ions in the
306 window, effectively removing background noise. Data centroiding is performed if necessary to
307 reduce data complexity. Ion images are extracted using mass bins of 0.01 m/z , and then spatially
308 correlated. A minimum of 5 shared pixels with non-zero intensities between two ion images is
309 required to ensure meaningful correlations and mitigate the impact of sparse data. Pseudo
310 MS/MS spectra are generated by applying a predefined spatial correlation cutoff (e.g., 0.85),
311 followed by deduplication to remove redundant spectra. Both Numba⁵² acceleration and parallel
312 processing are employed for computation efficiency enhancement.

313

314 ***Reverse spectral search***

315 Reverse spectral search is an asymmetric matching process, where one spectrum is
316 treated as template (**T**) and the other as query (**Q**). All the peaks in the template spectrum and
317 aligned peaks in the query spectrum are involved in matching score calculation, shown as follows.

318
$$reverse\ cosine = \frac{Q_{aligned} \cdot T_{aligned}}{\|Q_{aligned}\| \|T_{all}\|}$$

319 Considering that pseudo MS/MS spectra are generated from low-energy fragmentation
320 scan modes, and that most public reference MS/MS are acquired under medium to high collision
321 energies, we propose a mass-dependent approach to scale peak intensities for reference MS/MS
322 spectra. This method aims to simulate the pattern observed in low-energy MS/MS, where
323 fragment ions with m/z values closer to the precursor m/z exhibit higher abundance, while those
324 further from the precursor m/z show lower abundance. For a reference MS/MS, we have

325

$$\alpha_i = \frac{m/z_i}{m/z_{precursor}}$$

326

$$I_{i,new} = e^{k\alpha_i} \cdot I_i$$

327 where α_i is the m/z ratio of fragment i to the precursor; I_i is the original intensity; $I_{i,new}$ is the
328 scaled peak intensity; k is the scaling factor where we set it as 8 throughout this paper. Square
329 root transformation is then applied on both reference and pseudo MS/MS spectra. Each pseudo
330 MS/MS is searched against the reference MS/MS library in a precursor-tolerant manner, where
331 we ask that the precursor ions of matching hits should be in the m/z values of the query pseudo
332 MS/MS spectrum. For each pseudo MS/MS, we reserve the top 1 hit for each unique precursor
333 m/z value among all annotations. Each annotation is then linked to a single metabolic feature in
334 the metabolic feature table using the retention time of the pseudo MS/MS and the precursor mass
335 of the annotation.

336

337 To speed up the process of library search, we modified the flash entropy framework³¹
338 specifically for reverse cosine search, which outputs reverse cosine score, matched peak number
339 and spectral usage (sum intensities of matched peaks over total intensities in the query spectrum).
340 The following cutoffs were used for LC-MS MS1 data annotation: minimum score, 0.7; minimum
341 matched peaks, 4; minimum spectral usage, 0.20. For MS imaging data, we used: minimum
342 score, 0.7; minimum matched peaks, 4; minimum spectral usage, 0.05.

343

344 The reference MS/MS library needs to be preprocessed and indexed before use. We
345 provided the indexed version of GNPS MS/MS library (downloaded on July 17, 2024) as well as
346 the code to index an MS/MS library on GitHub (https://github.com/Philipbear/ms1_id).

347

348 **Preparation of chemical standards**

349 For the bile acid pool, a stock solution of 10 mM of glycocholic acid (GCA),
350 glycochenodeoxycholic acid (GCDCA), taurodeoxycholic acid (TDCA), taurocholic acid (TCA),
351 tauroolithocholic acid (TLCA), and tauro-3 α hydroxy-12ketocholeic acid was prepared. All
352 bile acids were diluted to 10 μ M into a single 2 mL LC-MS glass vial (Thermo Fisher) to create a
353 pooled sample. For the drug pool, a stock solution of 10 mM of sertraline, venlafaxine, ritonavir,
354 darunavir, losartan, quetiapine, sulfasalazine, and abacavir was prepared. All drugs were diluted
355 to 10 μ M into a single 2 mL LC-MS glass vial (Thermo Fisher) to create a pooled sample.

356

357 **Preparation of NIST reference materials**

358 NIST fecal reference materials (two vegan tubes and two omnivore tubes) were subjected
359 to a biphasic extraction⁵³ to remove lipids and retain the metabolite fraction. One mL of NIST fecal
360 material was transferred to a 2 mL Eppendorf tube and dried overnight in a CentriVap. Dry
361 materials were resuspended with 325 μ L of cold MeOH (LC-MS grade, Thermo Fisher), vortex for
362 10 s, and sonicated for 5 min before adding 1083 μ L of cold MTBE. Samples were vortexed for

363 10 s and sonicated for 2 min followed by 1 h incubation at 4 °C. To induce phase separation, 271
364 μL H_2O (LC-MS grade, Thermo Fisher) was added to the samples and centrifuge at 10,000 x g
365 for 10 min. The upper phase was removed and 1084 μL of MeOH was added, followed by an
366 overnight incubation at -20 °C. Samples were centrifuged at 15,000 x g for 10 min. An equal
367 amount (50 μL) of the fecal NIST materials were combined to generate a pooled NIST reference
368 fecal sample. All samples were dried in a CentriVap and stored at -80 °C until resuspension. NIST
369 reference fecal materials were resuspended in 200 μL of 50% MeOH/ H_2O with sulfadimethoxine
370 as internal standard before LC-MS analysis.

371

372 ***LC-MS analysis***

373 The chromatographic separation was done on a reverse phase polar C18 (Kinetex Polar
374 C18, 100 mm x 2.1 mm, 2.6 μm , 100 angstrom pore size with the matching guard column,
375 Phenomenex) using a Vanquish UHPLC coupled to an Orbitrap mass spectrometer (Thermo
376 Fisher Scientific). Five microliters of samples were injected into the mobile phase, which is
377 composed of solvent A (H_2O with 0.1% formic acid) and solvent B (ACN with 0.1% formic acid)
378 with the column compartment kept at 40 °C. Samples were eluted at a flow rate of 0.5 mL/min
379 using the following gradient: 0 min, 5% B; 1.1 min, 5% B; 7.5 min, 40% B; 8.5 min, 99% B; 9.5
380 min, 99% B; 10 min, 5% B; 10.5 min, 5% B; 10.75 min, 99% B; 11.25 min, 99% B; 11.5 min, 5%
381 B; 12.5 min, 5% B. Data were acquired using DDA mode or full-scan mode in electrospray positive
382 ionization mode.

383

384 For DDA mode, the parameters were set as: sheath gas flow 53 L/min, aux gas flow rate
385 14 L/min, sweep gas flow 3 L/min, spray voltage 3.5 kV, inlet capillary 269 °C, aux gas heater
386 438 °C, S-lens RF level 50.0. MS scan range was set as 100-1000 m/z with mass resolution of
387 35,000 at m/z 200. Automatic gain control (AGC) target was set to 1E6 with a maximum injection
388 time of 100 ms. Up to 5 MS/MS spectra per MS1 were collected per cycle with mass resolution
389 17,500 at m/z 200, maximum injection time of 150 ms with an AGC target of 5E5. Isolation window
390 was set to 1 m/z and the isolation offset at 0 m/z . Stepwisely normalized collision energies were
391 set at 25 eV, 40 eV, and 60 eV. The apex trigger was set to 2-15 s and a dynamic exclusion of 5
392 s. Isotopes were excluded from the analysis.

393

394 For full-scan mode, the parameters were set as: sheath gas flow 53 L/min, aux gas flow
395 rate 14 L/min, sweep gas flow 3 L/min, spray voltage 3.5 kV, inlet capillary 269 °C, and aux gas
396 heater 438 °C. MS scan range was set as 100-1000 m/z with mass resolution of 70,000 at m/z
397 200. AGC target was set to 1E6 with maximum injection time as 150 ms. Data in full-scan mode
398 were acquired using different isCID energies: 0 eV, 10 eV, and 20 eV.

399

400 ***Statistical analysis***

401 positive mode). For each metabolic feature, we cleaned its pseudo MS/MS by removing
402 all ions larger. For the IBD dataset, missing values were filled using the minimum of 5E5 and 10%
403 of the minimum intensity for each feature. Outlier removal was conducted using the interquartile
404 range (IQR) method. Data points below $Q1 - 1.5 \times IQR$ or above $Q3 + 1.5 \times IQR$ were removed
405 from each group (non-IBD, CD or UC). Two-side Mann-Whitney U tests were performed. For t-
406 SNE visualization, intensity values were subjected to log transformation and feature-wise z-
407 normalization.

408

409 ***Molecular networking***

410 A molecular network was constructed using the pseudo MS/MS spectra from the IBD
411 dataset (HILIC positive mode). For each metabolic feature, we cleaned its pseudo MS/MS by
412 removing all ions larger than the feature m/z . Then, an MGF file for all cleaned pseudo MS/MS
413 spectra was prepared. The MGF file was uploaded onto the GNPS2 platform, where a classical
414 molecular networking workflow (version 2024.09.20) was completed. A minimum of modified
415 cosine of 0.8 and matched peaks of 6 are required to build an edge in the network construction.
416 The job is available at <https://gnps2.org/status?task=670aa34a07544a5cbbd1f1d40605f50f>.

417

418

419 **Data availability**

420 All the source data used in this study are publicly accessible. For LC-MS data, pooled
421 chemical standards are available at GNPS/MassIVE repository with accession number
422 [MSV000095789](https://massive.ucsf.edu/MSV000095789); NIST human feces are available at GNPS/MassIVE repository with accession
423 number [MSV000095787](https://massive.ucsf.edu/MSV000095787); the IBD dataset is available at Metabolomics Workbench with project
424 number [PR000639](https://www.ebi.ac.uk/metabolomics/projects/PR000639). For MS imaging data, the mouse brain data are available at MetaboLights
425 repository under code [MTBLS313](https://www.ebi.ac.uk/metabolomics/projects/MTBLS313); the mouse body dataset is available at METASPACE platform
426 with ID [2022-07-08_20h45m00s](https://metaspace.org/2022-07-08_20h45m00s); the hepatocytes data are available at METASPACE platform
427 with project ID [Rappez_2021_SpaceM](https://metaspace.org/Rappez_2021_SpaceM).

428

429

430 **Code availability**

431 Source codes are available at GitHub (https://github.com/Philipbear/ms1_id) and Zenodo
432 (<https://zenodo.org/records/13864878>) under the Apache-2.0 license.

433

434

435 **Acknowledgements**

436 This work was supported by BBSRC/NSF award 2152526 and National Institute of Health
437 Sciences U24DK133658. We acknowledge the NIST Complex Microbial Systems Group for
438 providing the NIST material ahead of its official release. Thanks to Theodore Alexandrov for
439 guiding us to the datasets in METASPACE.

440

441

442 **Disclosures**

443 P.C.D. is a scientific advisor and holds equity in Sirenas, Cybele, and bileOmix, and is a
444 Scientific Co-founder, and advisor and holds equity in Ometa, Arome, and Enveda with prior
445 approval by UC-San Diego.

446

447

448 **Author contributions**

449 P.C.D. and S.X. conceived the research project. S.X. developed the computational
450 algorithm and performed data analysis. V.C.L. collected the LC-MS data. Y.E. provided LC-MS
451 file summaries in public repositories. S.X. and P.C.D. drafted the manuscript. P.C.D. supervised
452 the project. All authors approved the manuscript.

453

454

455 **References**

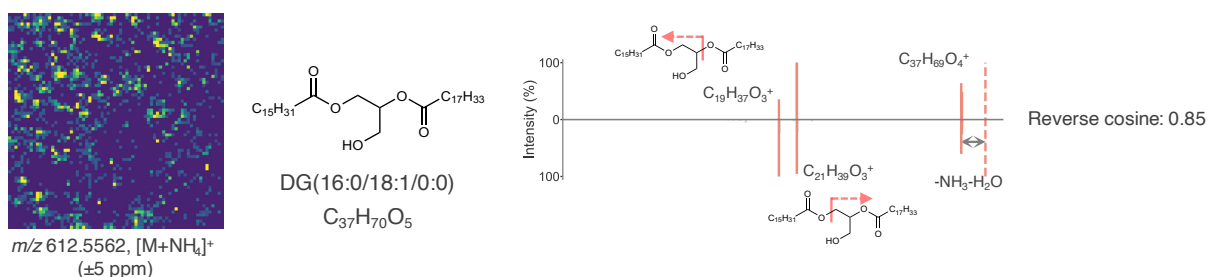
- 456 1. Sud, M. *et al.* Metabolomics Workbench: An international repository for metabolomics data
457 and metadata, metabolite standards, protocols, tutorials and training, and analysis tools.
458 *Nucleic Acids Res.* **44**, D463–D470 (2016).
- 459 2. Yurekten, O. *et al.* MetaboLights: open data repository for metabolomics. *Nucleic Acids Res.*
460 **52**, D640–D646 (2024).
- 461 3. Wang, M. *et al.* Sharing and community curation of mass spectrometry data with Global
462 Natural Products Social Molecular Networking. *Nat. Biotechnol.* **34**, 828–837 (2016).
- 463 4. Palmer, A. *et al.* FDR-controlled metabolite annotation for high-resolution imaging mass
464 spectrometry. *Nat. Methods* **14**, 57–60 (2017).
- 465 5. da Silva, R. R., Dorrestein, P. C. & Quinn, R. A. Illuminating the dark matter in metabolomics.
466 *Proc. Natl. Acad. Sci.* **112**, 12549–12550 (2015).
- 467 6. Blaženović, I. *et al.* Structure Annotation of All Mass Spectra in Untargeted Metabolomics.
468 *Anal. Chem.* **91**, 2155–2162 (2019).
- 469 7. Aksenov, A. A., da Silva, R., Knight, R., Lopes, N. P. & Dorrestein, P. C. Global chemical
470 analysis of biology by mass spectrometry. *Nat. Rev. Chem.* **1**, 1–20 (2017).
- 471 8. Bittremieux, W., Wang, M. & Dorrestein, P. C. The critical role that spectral libraries play in
472 capturing the metabolomics community knowledge. *Metabolomics* **18**, 94 (2022).
- 473 9. Dührkop, K. *et al.* SIRIUS 4: a rapid tool for turning tandem mass spectra into metabolite
474 structure information. *Nat. Methods* **16**, 299–302 (2019).
- 475 10. Xing, S., Shen, S., Xu, B., Li, X. & Huan, T. BUDDY: molecular formula discovery via
476 bottom-up MS/MS interrogation. *Nat. Methods* 1–10 (2023) doi:10.1038/s41592-023-01850-
477 x.
- 478 11. Gabelica, V. & Pauw, E. D. Internal energy and fragmentation of ions produced in

- 479 electrospray sources. *Mass Spectrom. Rev.* **24**, 566–587 (2005).
- 480 12. Abrankó, L., García-Reyes, J. F. & Molina-Díaz, A. In-source fragmentation and
481 accurate mass analysis of multiclass flavonoid conjugates by electrospray ionization time-of-
482 flight mass spectrometry. *J. Mass Spectrom.* **46**, 478–488 (2011).
- 483 13. Criscuolo, A., Zeller, M. & Fedorova, M. Evaluation of Lipid In-Source Fragmentation on
484 Different Orbitrap-based Mass Spectrometers. *J. Am. Soc. Mass Spectrom.* **31**, 463–466
485 (2020).
- 486 14. Xu, Y.-F., Lu, W. & Rabinowitz, J. D. Avoiding Misannotation of In-Source Fragmentation
487 Products as Cellular Metabolites in Liquid Chromatography–Mass Spectrometry-Based
488 Metabolomics. *Anal. Chem.* **87**, 2273–2281 (2015).
- 489 15. Giera, M., Aisporna, A., Uritboonthai, W. & Siuzdak, G. The hidden impact of in-source
490 fragmentation in metabolic and chemical mass spectrometry data interpretation. *Nat. Metab.*
491 1–2 (2024) doi:10.1038/s42255-024-01076-x.
- 492 16. Domingo-Almenara, X. *et al.* Autonomous METLIN-Guided In-source Fragment
493 Annotation for Untargeted Metabolomics. *Anal. Chem.* **91**, 3246–3253 (2019).
- 494 17. Baygi, S. F., Kumar, Y. & Barupal, D. K. IDSL.CSA: Composite Spectra Analysis for
495 Chemical Annotation of Untargeted Metabolomics Datasets. *Anal. Chem.* **95**, 9480–9487
496 (2023).
- 497 18. Xue, J. *et al.* EISA-EXPOSOME: One Highly Sensitive and Autonomous Exposomic
498 Platform with Enhanced in-Source Fragmentation/Annotation. *Anal. Chem.* **95**, 17228–17237
499 (2023).
- 500 19. Wang, X.-C. *et al.* AntDAS-DDA: A New Platform for Data-Dependent Acquisition Mode-
501 Based Untargeted Metabolomic Profiling Analysis with Advantage of Recognizing Insource
502 Fragment Ions to Improve Compound Identification. *Anal. Chem.* **95**, 638–649 (2023).
- 503 20. Wasito, H., Causon, T. & Hann, S. Alternating in-source fragmentation with single-stage
504 high-resolution mass spectrometry with high annotation confidence in non-targeted
505 metabolomics. *Talanta* **236**, 122828 (2022).
- 506 21. Broeckling, C. D., Afsar, F. A., Neumann, S., Ben-Hur, A. & Prenni, J. E. RAMClust: A
507 Novel Feature Clustering Method Enables Spectral-Matching-Based Annotation for
508 Metabolomics Data. *Anal. Chem.* **86**, 6812–6817 (2014).
- 509 22. Graça, G. *et al.* Automated Annotation of Untargeted All-Ion Fragmentation LC–MS
510 Metabolomics Data with MetaboAnnotatoR. *Anal. Chem.* **94**, 3446–3455 (2022).
- 511 23. Kachman, M. *et al.* Deep annotation of untargeted LC-MS metabolomics data with
512 Binner. *Bioinformatics* **36**, 1801–1806 (2020).
- 513 24. Li, Y. *et al.* Spectral entropy outperforms MS/MS dot product similarity for small-
514 molecule compound identification. *Nat. Methods* **18**, 1524–1531 (2021).
- 515 25. Aksenov, A. A. *et al.* Auto-deconvolution and molecular networking of gas
516 chromatography–mass spectrometry data. *Nat. Biotechnol.* **39**, 169–173 (2021).
- 517 26. Xue, J. *et al.* Enhanced in-Source Fragmentation Annotation Enables Novel Data

- 518 Independent Acquisition and Autonomous METLIN Molecular Identification. *Anal. Chem.* **92**,
519 6051–6059 (2020).
- 520 27. Wishart, D. S. *et al.* HMDB 5.0: the Human Metabolome Database for 2022. *Nucleic*
521 *Acids Res.* **50**, D622–D631 (2022).
- 522 28. Sud, M. *et al.* LMSD: LIPID MAPS structure database. *Nucleic Acids Res.* **35**, D527–
523 D532 (2007).
- 524 29. Hastings, J. *et al.* ChEBI in 2016: Improved services and an expanding collection of
525 metabolites. *Nucleic Acids Res.* **44**, D1214–D1219 (2016).
- 526 30. Schymanski, E. L. *et al.* Identifying Small Molecules via High Resolution Mass
527 Spectrometry: Communicating Confidence. *Environ. Sci. Technol.* **48**, 2097–2098 (2014).
- 528 31. Li, Y. & Fiehn, O. Flash entropy search to query all mass spectral libraries in real time.
529 *Nat. Methods* 1–4 (2023) doi:10.1038/s41592-023-02012-9.
- 530 32. Tsugawa, H. *et al.* MS-DIAL: data-independent MS/MS deconvolution for comprehensive
531 metabolome analysis. *Nat. Methods* **12**, 523–526 (2015).
- 532 33. Kuhl, C., Tautenhahn, R., Böttcher, C., Larson, T. R. & Neumann, S. CAMERA: An
533 Integrated Strategy for Compound Spectra Extraction and Annotation of Liquid
534 Chromatography/Mass Spectrometry Data Sets. *Anal. Chem.* **84**, 283–289 (2012).
- 535 34. Mahieu, N. G. & Patti, G. J. Systems-Level Annotation of a Metabolomics Data Set
536 Reduces 25 000 Features to Fewer than 1000 Unique Metabolites. *Anal. Chem.* **89**, 10397–
537 10406 (2017).
- 538 35. Nash, W. J., Ngere, J. B., Najdekr, L. & Dunn, W. B. Characterization of Electrospray
539 Ionization Complexity in Untargeted Metabolomic Studies. *Anal. Chem.* (2024)
540 doi:10.1021/acs.analchem.4c00966.
- 541 36. Djoumbou Feunang, Y. *et al.* ClassyFire: automated chemical classification with a
542 comprehensive, computable taxonomy. *J. Cheminformatics* **8**, 61 (2016).
- 543 37. Lloyd-Price, J. *et al.* Multi-omics of the gut microbial ecosystem in inflammatory bowel
544 diseases. *Nature* **569**, 655–662 (2019).
- 545 38. Kim, H. W. *et al.* NPClassifier: A Deep Neural Network-Based Structural Classification
546 Tool for Natural Products. *J. Nat. Prod.* **84**, 2795–2807 (2021).
- 547 39. Lavelle, A. & Sokol, H. Gut microbiota-derived metabolites as key actors in inflammatory
548 bowel disease. *Nat. Rev. Gastroenterol. Hepatol.* **17**, 223–237 (2020).
- 549 40. Chang, F.-Y. *et al.* Gut-inhabiting Clostridia build human GPCR ligands by conjugating
550 neurotransmitters with diet- and human-derived fatty acids. *Nat. Microbiol.* **6**, 792–805
551 (2021).
- 552 41. Cohen, L. J. *et al.* Commensal bacteria make GPCR ligands that mimic human signalling
553 molecules. *Nature* **549**, 48–53 (2017).
- 554 42. Rappez, L. *et al.* SpaceM reveals metabolic states of single cells. *Nat. Methods* **18**, 799–
555 805 (2021).
- 556 43. Derave, W., Everaert, I., Beeckman, S. & Baguet, A. Muscle Carnosine Metabolism and

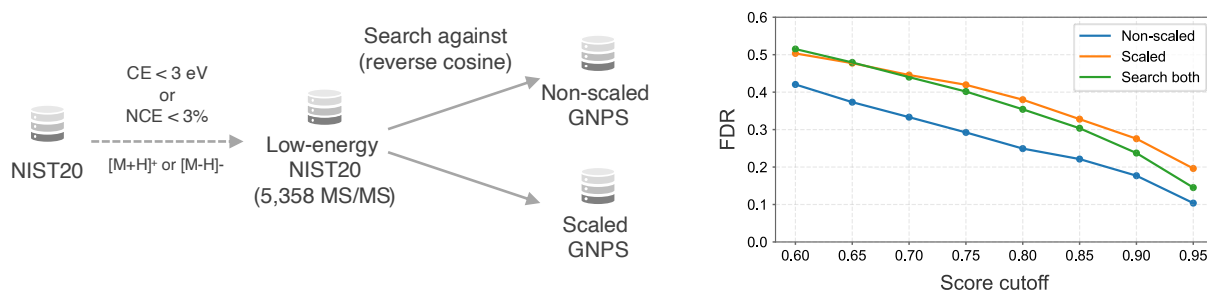
- 557 β -Alanine Supplementation in Relation to Exercise and Training. *Sports Med.* **40**, 247–263
558 (2010).
- 559 44. Boldyrev, A. A., Aldini, G. & Derave, W. Physiology and Pathophysiology of Carnosine.
560 *Physiol. Rev.* **93**, 1803–1845 (2013).
- 561 45. Guney, Y. *et al.* Carnosine may reduce lung injury caused by radiation therapy. *Med.*
562 *Hypotheses* **66**, 957–959 (2006).
- 563 46. Wang, M. *et al.* Mass spectrometry searches using MASST. *Nat. Biotechnol.* **38**, 23–26
564 (2020).
- 565 47. Zuffa, S. *et al.* microbeMASST: a taxonomically informed mass spectrometry search tool
566 for microbial metabolomics data. *Nat. Microbiol.* 1–10 (2024) doi:10.1038/s41564-023-01575-
567 9.
- 568 48. Gentry, E. C. *et al.* Reverse metabolomics for the discovery of chemical structures from
569 humans. *Nature* 1–8 (2023) doi:10.1038/s41586-023-06906-8.
- 570 49. Abiead, Y. E. *et al.* Enabling pan-repository reanalysis for big data science of public
571 metabolomics data. Preprint at <https://doi.org/10.26434/chemrxiv-2024-jt46s> (2024).
- 572 50. Bittremieux, W. *et al.* Open access repository-scale propagated nearest neighbor
573 suspect spectral library for untargeted metabolomics. *Nat. Commun.* **14**, 8488 (2023).
- 574 51. Yu, H. <https://github.com/huaxuyu/masscube>. (2024).
- 575 52. Lam, S. K., Pitrou, A. & Seibert, S. Numba: a LLVM-based Python JIT compiler. in
576 *Proceedings of the Second Workshop on the LLVM Compiler Infrastructure in HPC* 1–6
577 (Association for Computing Machinery, New York, NY, USA, 2015).
578 doi:10.1145/2833157.2833162.
- 579 53. Coman, C. *et al.* Simultaneous Metabolite, Protein, Lipid Extraction (SIMPLEX): A
580 Combinatorial Multimolecular Omics Approach for Systems Biology*. *Mol. Cell. Proteomics*
581 **15**, 1435–1466 (2016).
582

583
584



585
586 **Extended Data Fig. 1** | DG(16:0/18:1/0:0) was annotated in hepatocytes MS imaging data, with a reverse cosine score of 0.85. Both
587 acyl chains could be annotated in the MS/MS spectrum. The ion image was created using a mass tolerance of 5 ppm.

588
589
590
591



592
593 **Extended Data Fig. 2** | Peak scaling tends to increase the false discovery rate (FDR) during spectral matching. We selected low-
594 energy reference MS/MS spectra from the NIST20 MS/MS library using collision energy (CE) < 3 eV or normalized collision energy
595 (NCE) < 3%. These spectra were then searched against the non-scaled and scaled GNPS library. With a minimum of 4 matched
596 peaks (peaks other than the precursor ion), FDR results of different score cutoffs were shown in the line plot. We expect future
597 optimization of scaling can further improve the annotation confidence.

598
599
600



601
602 **Extended Data Fig. 3** | A proof-of-principle of MS1-based MASST. We searched the pseudo MS/MS spectrum of phenylalanine-C3:0
603 from NIST human feces data against the pseudo MS/MS spectra pool from the IBD dataset, and it returned a match of cosine 0.90
604 with 6 matched peaks. The returned hit was also annotated as phenylalanine-C3:0 in the IBD dataset.

605