

HERMES: a molecular formula-oriented method to target the metabolome

Roger Giné¹, Jordi Capellades^{1,2}, Josep M. Badia^{1,2}, Dennis Vughs³, Michaela Schwaiger-Haber^{4,5}, Theodore Alexandrov^{6,7,8}, Maria Vinaixa^{1,2}, Andrea M. Brunner³, Gary J. Patti^{4,5}, Oscar Yanes^{*1,2}

1. Universitat Rovira i Virgili, Department of Electronic Engineering & IISPV, Tarragona, Spain.
2. CIBER de Diabetes y Enfermedades Metabólicas Asociadas (CIBERDEM), Instituto de Salud Carlos III, Madrid, Spain.
3. KWR Water Research Institute, Nieuwegein, The Netherlands.
4. Department of Chemistry, Washington University, St. Louis, MO, USA.
5. Department of Medicine, Washington University, St. Louis, MO, USA.
6. Structural and Computational Biology Unit, European Molecular Biology Laboratory, Heidelberg, Germany.
7. Molecular Medicine Partnership Unit, European Molecular Biology Laboratory, Heidelberg, Germany.
8. Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California San Diego, La Jolla, California, USA.

*Corresponding author:

Oscar Yanes, PhD

Department of Electronic Engineering

Universitat Rovira i Virgili

Avinguda Països Catalans, 26, 43007 Tarragona, Spain

Phone: +34 977759397

Email: oscar.yanes@urv.cat

35 **Comprehensive metabolome analyses are essential for biomedical, environmental and**
36 **biotechnological research. However, current MS1 and MS2-based acquisition and data**
37 **analysis strategies in untargeted metabolomics result in low identification rates of**
38 **metabolites. Here we present HERMES, a molecular formula-oriented and peak**
39 **detection-free method that uses raw LC/MS1 information to optimize MS2 acquisition.**
40 **Investigating environmental water, *E. coli*, and human plasma extracts with HERMES,**
41 **we achieved an increased biological specificity of MS2 scans, leading to improved**
42 **mass spectral similarity scoring and identification rates when compared to a state-of-**
43 **the-art data-dependent acquisition (DDA) approach. Thus, HERMES improves**
44 **sensitivity, selectivity and annotation of metabolites. HERMES is available as an R**
45 **package with a user-friendly graphical interface for data analysis and visualization.**

46

47 **Introduction**

48 A single LC/MS-based metabolomic experiment generates millions of three-
49 dimensional (m/z , retention time, intensity) data points that can be annotated and quantified
50 into thousands of metabolite features. However, most features are either redundant ions
51 caused by ionization-related phenomena such as cation/anion adduction, multimerization and
52 in-source fragmentation, or unknown contaminants and artifacts^{1,2}. Moreover, conventional
53 untargeted metabolomic experiments lead to highly heterogeneous chromatographic peak
54 shapes, which negatively affect the performance of peak detection³ and grouping/annotation
55 algorithms in MS1 mode⁴. These characteristics of MS1 data, in turn, negatively impact MS2
56 acquisition methods used for metabolite identification. In data-dependent acquisition (DDA)
57 mode, MS2 spectra are automatically collected for precursor ions that exceed a predefined
58 intensity threshold. The selection of precursor ions is a stochastic event suffering from low
59 analytical reproducibility and favouring the selection of the most abundant, but not necessarily
60 biologically relevant, ions. In data-independent acquisition (DIA) methods, multiple precursor
61 ions, including redundant and biologically irrelevant ions, are simultaneously fragmented,
62 often generating a series of complex convoluted MS2 spectra. Despite the emergence of new
63 software to reconstruct the link between precursors and their fragments through mass spectral
64 deconvolution^{5,6}, MS2 spectral quality and matching scores to reference spectra are generally
65 poorer in DIA compared to DDA⁷.

66

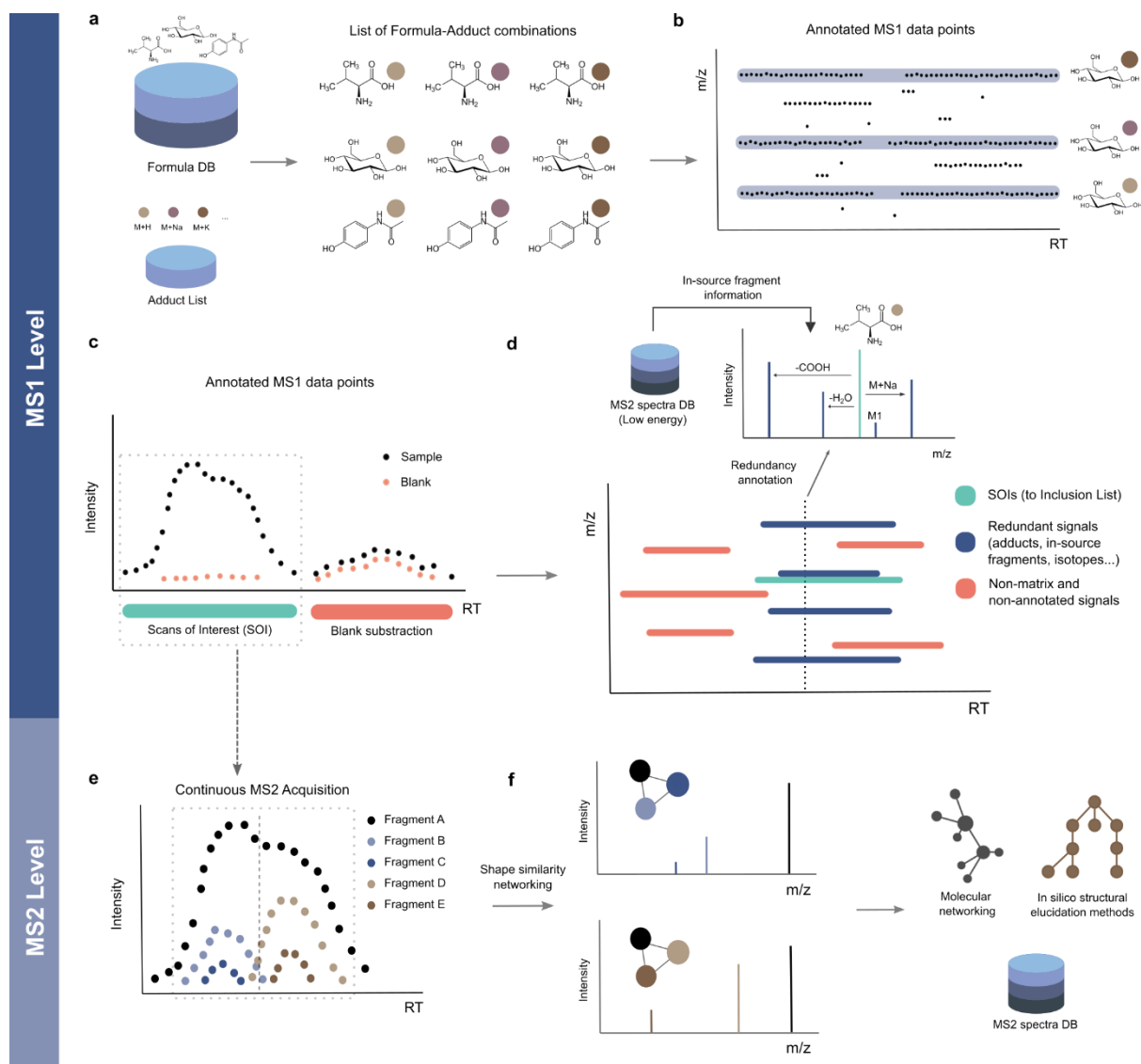
67

68

69 Results

70 Here we present HERMES, a novel experimental method and computational tool in
71 untargeted metabolomics that improves the selectivity and sensitivity for comprehensive
72 metabolite profiling in MS1, and identification in MS2. HERMES replaces the conventional
73 untargeted metabolomic workflow that detects and annotates peaks^{8,9}, for an inverse
74 approach that directly interrogates raw LC/MS1 data points (within scans) by using a
75 comprehensive list of unique molecular formulas selected by the user. These are retrieved
76 from large compound-centric databases (e.g., HMDB, ChEBI, NORMAN)¹⁰⁻¹², genome-scale
77 metabolic models, or specific metabolic pathways. Each molecular formula generates multiple
78 'ionic formulas' by adding or subtracting atoms from common adduct ions (Fig. 1). The
79 resulting ionic formulas (on the order of 10^4 - 10^5 from a database such as HMDB) are searched
80 against millions of data points in an LC/MS1 experiment. HERMES calculates the theoretical
81 isotopic pattern of each ionic formula based on a predefined experimental mass resolution
82 value (Suppl. Fig. 1). The number of collisions between monoisotopic ionic formulas vary
83 according to the experimental mass error (i.e., the smaller the error, the larger the percentage
84 of non-overlapping ionic formulas; Suppl. Fig. 2). An LC/MS1 data point contains m/z and
85 intensity information in a wide mass range (e.g., m/z 80 to 1,000) for a given instant of time
86 (Suppl. Fig 3). HERMES solves the limitations of peak detection by finding a series of scans,
87 named SOI (Scans Of Interest), which are defined as clusters of data points that match an
88 ionic formula, are concentrated within a short period of time and contain a minimum amount
89 of structure determined by a 1D version of the ρ_{chaos} score from Palmer et al.¹³ (see Online
90 Methods). SOI shapes do not necessarily fit a Gaussian-like function, as assumed in basic
91 chromatography theory, making the process independent of the heterogeneous peak shapes
92 commonly observed in LC/MS1 experiments from complex mixtures. SOIs are then filtered in
93 three steps: (i) blank subtraction from the sample based on an artificial neural network (Suppl.
94 Fig. 4a), (ii) adduct and isotopologue grouping according to the similarity of their elution
95 profiles (Suppl. Fig. 4b), and (iii) in-source fragment (ISF) annotation by using publicly
96 available low-energy MS2 data (Suppl. Fig. 4c) extending on Domingo-Almenara et al.¹⁴.
97 Finally, users can prioritize the SOIs that will constitute the inclusion list for targeted MS2
98 acquisition based on the following criteria: type and number of adducts, minimum intensity,
99 isotopic fidelity, and a maximum number of overlapped precursors at any time range, which
100 together determine the total number of MS2 runs. According to the MS2 acquisition settings,
101 each entry in the inclusion list may be associated with one or multiple MS2 scans: if there are
102 more than five continuous scans, HERMES provides an optional deconvolution step (adapted
103 from CliqueMS¹⁵) that resolves partially co-eluting isomeric compounds (Suppl. Fig. 5); if there
104 are fewer scans, HERMES selects the most intense scan. The resulting curated MS2 spectra

105 can either be identified within HERMES or exported as .mzML, .msp, or .mgf files to be used
 106 in other identification software such as NIST MS Search, SIRIUS^{16,17}, or GNPS¹⁸.



107

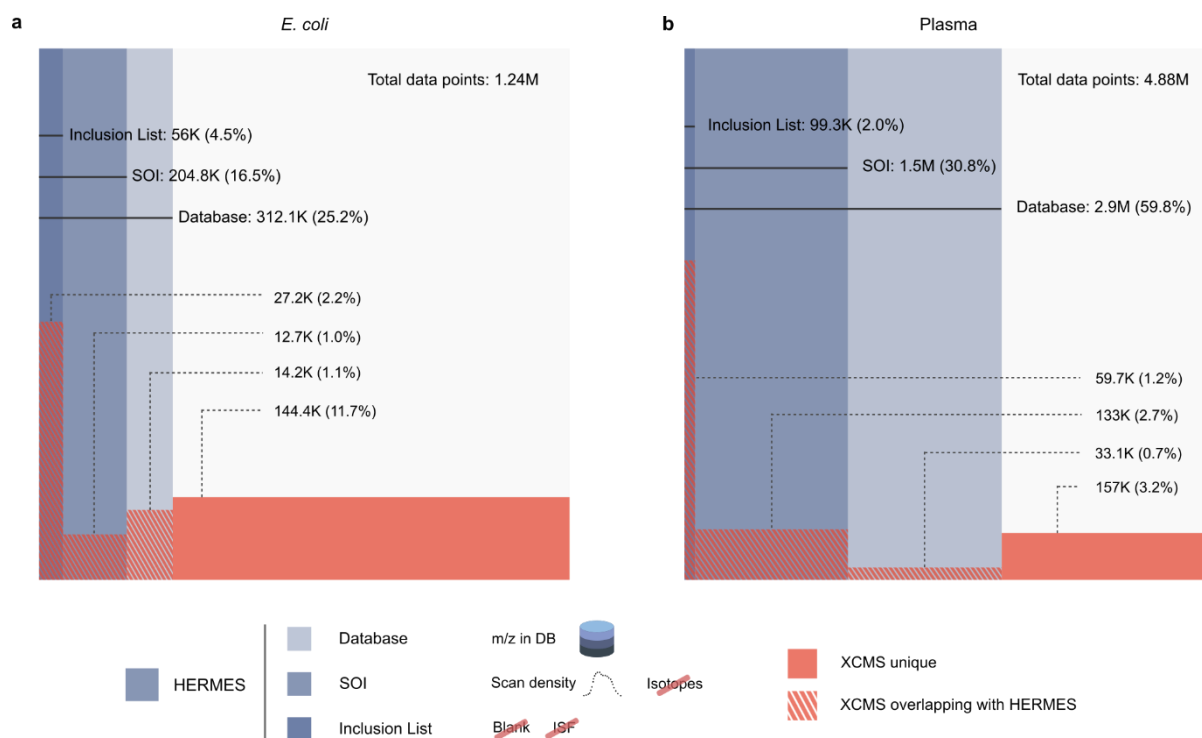
108 **Figure 1. The HERMES workflow.** (a) A context-specific database of molecular formulas and
 109 MS adducts generates a list of ionic formulas. (b) LC/MS1 data points are interrogated against all m/z
 110 ions corresponding to the ionic formulas and their isotopes. (c) Points with the same m/z annotation are
 111 grouped by density into retention time (RT) intervals called Scans of Interest (SOI). SOIs with similar
 112 shape and intensity in a blank sample are removed. (d) SOIs corresponding to different adducts of the
 113 same formula are grouped by their chromatographic elution profile. Similarly, in-source fragments are
 114 annotated based on low intensity MS2 spectra of molecules with the same formula. The result is an
 115 inclusion list (IL) of sample-specific and non-redundant precursor ions that will be monitored in a
 116 posterior MS2 experiment. (e) The IL entries are acquired continuously along the defined RT interval
 117 and HERMES groups the resulting fragment elution profiles. (f) This results in deconvoluted spectra of
 118 partially co-eluting isomeric compounds that can be queried against an MS2 database or exported to
 119 be used in alternative identification workflows.

120 HERMES is available as an R package (RHermes) and comes with an R Graphical
121 User Interface (GUI) to allow data analysis, tracking of compound annotations, and
122 visualization (Suppl. Fig. 6). RHermes accepts both CSV and XLS/XLSX files as valid
123 molecular formula lists and can extract formulas from selected KEGG pathways for a given
124 organism. RHermes is designed to work on ≥ 8 GB RAM computers, however the size of the
125 database of molecular formulas determines CPU usage and RAM memory. The running time,
126 including blank subtraction and inclusion list generation, goes from 1 to 20 minutes on a six-
127 core, 2.9 GHz CPU (Suppl. Table 1).

128 HERMES has been validated by using three (bio)chemically relevant samples of
129 increasing complexity: (i) water collected from a canal in Nieuwegein (Netherlands), (ii) *E.coli*,
130 and (iii) human plasma extracts. The canal water was spiked with 86 common environmental
131 contaminants at 1 $\mu\text{g/L}$ (Suppl. Table 2) and analysed by RP/LC (C18) coupled to an Orbitrap
132 in positive (pos) and negative (neg) ionization mode operating at 120,000 resolution. HERMES
133 detected and annotated all spiked compounds at the MS1 level using 118,820 (pos) and
134 46,809 (neg) ionic formulas calculated from 24,696 unique molecular formulas in the
135 NORMAN database. Certain ionic formula collisions, particularly those involving Cl, Br, S, or
136 K, were automatically resolved by matching experimental isotopic patterns to the expected
137 ones. This is the case, for example, of the $[M+H]^+$ ion of chloridazon and the $[M+K]^+$ ion of 2-
138 amino-alpha-carboline, which overlapped at 0.27 ppm (Suppl. Fig. 7). In-source fragments
139 that could be wrongly associated with ionic formulas were also annotated by using low-energy
140 MS2 spectra when available. The output was a curated inclusion list of 474 (pos) and 129
141 (neg) selective entries for targeted MS2 (Suppl. Fig. 8).

142 Next, a reference *E. coli* cell extract (Cambridge Isotope Laboratories) was analysed
143 by HILIC coupled to an Orbitrap in positive and negative ionization mode. LC/MS1 data were
144 analysed by HERMES by using 12,010 (pos) and 4,876 (neg) ionic formulas calculated from
145 2,463 unique molecular formulas obtained from the *Escherichia coli* Metabolome Database
146 (ECMDB) and KEGG database. Interestingly, HERMES annotated ionic formulas for 25%
147 (pos) and 22% (neg) of all data points acquired by the mass spectrometer (Fig. 2a and Suppl.
148 Fig. 9a). In comparison with XCMS, a commonly used open-source LC/MS1 processing data
149 tool in untargeted metabolomics^{9,19}, 16% of all acquired data points were associated with an
150 XCMS peak, 4.3% of data points in XCMS peaks matched an ionic formula from ECMDB and
151 KEGG database, and only 2.2% of data points in XCMS peaks were represented in the final
152 SOI list after blank subtraction, isotopic fidelity, and ISF removal. Consequently, the overlap
153 in isotopes and, more particularly, frequent adduct annotations between HERMES and well-
154 established tools for annotating LC/MS1 data^{15,20} was low: $\sim 80\%$ and $< 20\%$, respectively

155 (Suppl. Fig. 10a,b), showing the limitations of the peak detection strategies that have been
 156 used to date in untargeted metabolomics.



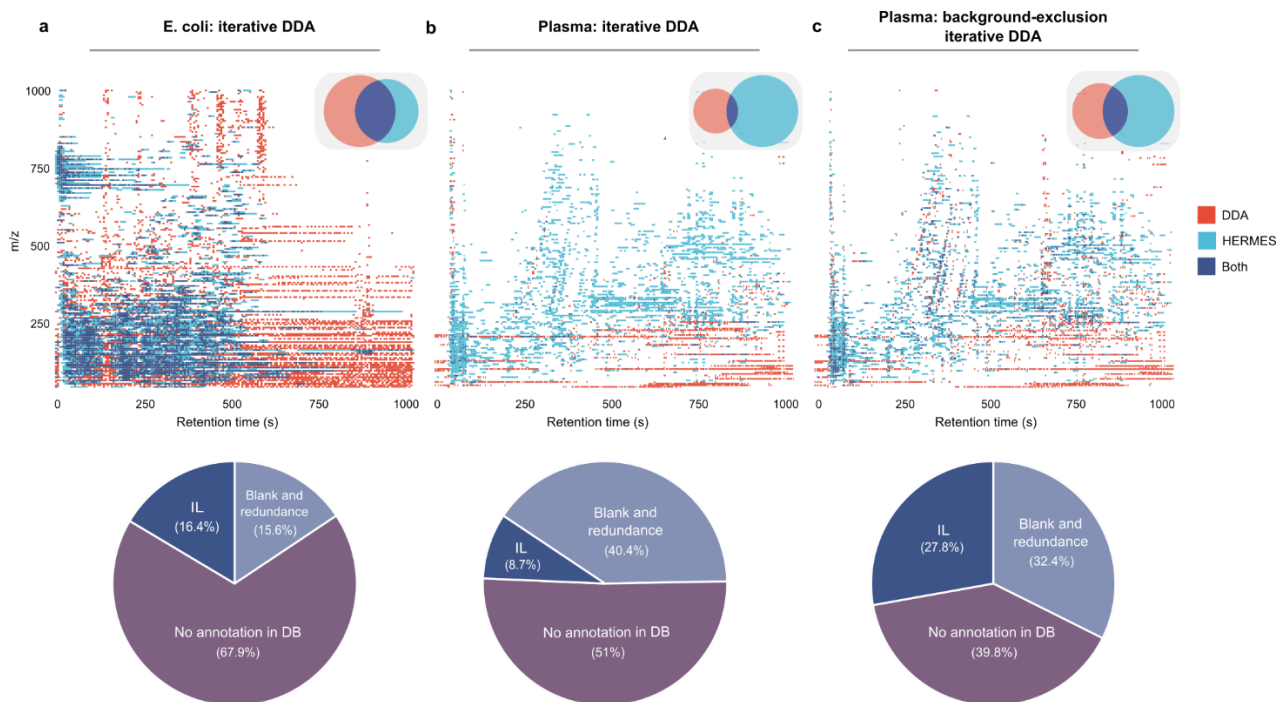
157

158 **Figure 2. Venn-like diagram of the distribution of LC/MS1 data points in different steps of the**
 159 **HERMES workflow and XCMS peak-associated points.** a) *E. coli* extract. b) Plasma extract.
 160 Database: Refers to all data points whose *m/z* matches with any *m/z* calculated from the ionic formula
 161 database (including isotopes). SOI: monoisotopic (M0)-annotated data points that are in Database and
 162 are also present in a SOI list that does not include blank subtraction nor any filtering. Inclusion List: data
 163 points present in Database and SOI kept through the blank subtraction, isotopic filter and ISF removal
 164 steps. Percentages refer to the total number of LC/MS1 data points. Positive ionization mode. On
 165 average, ~46% of data points in the inclusion list could not be annotated as a peak by XCMS.
 166 Conversely, ~86% scans annotated as a peak by XCMS could either not be matched to an ionic formula,
 167 were not specific of the sample or were associated with redundant signals.

168

169 The HERMES output was 2,058 (pos) and 1,081 (neg) SOIs that led to a curated
 170 inclusion list of 1,251 and 661 entries for targeted MS2, respectively. The *E. coli* extract was
 171 also analysed by iterative DDA under identical analytical conditions. Remarkably, 68% of DDA
 172 scans could not be annotated as the monoisotopic signal by any ionic formula from ECMDDB
 173 and KEGG database (Fig. 3a), which indicates their exogenous or artefactual origin. After
 174 filtering out DDA precursor ions that were classified as SOIs in the blank sample, redundant
 175 adducts, and ISF by HERMES; only 16% of the DDA scans matched with any monoisotopic

176 ionic formula in the inclusion list. In addition, HERMES included 591 inclusion list entries (47%
177 of the total) that were not triggered by DDA.



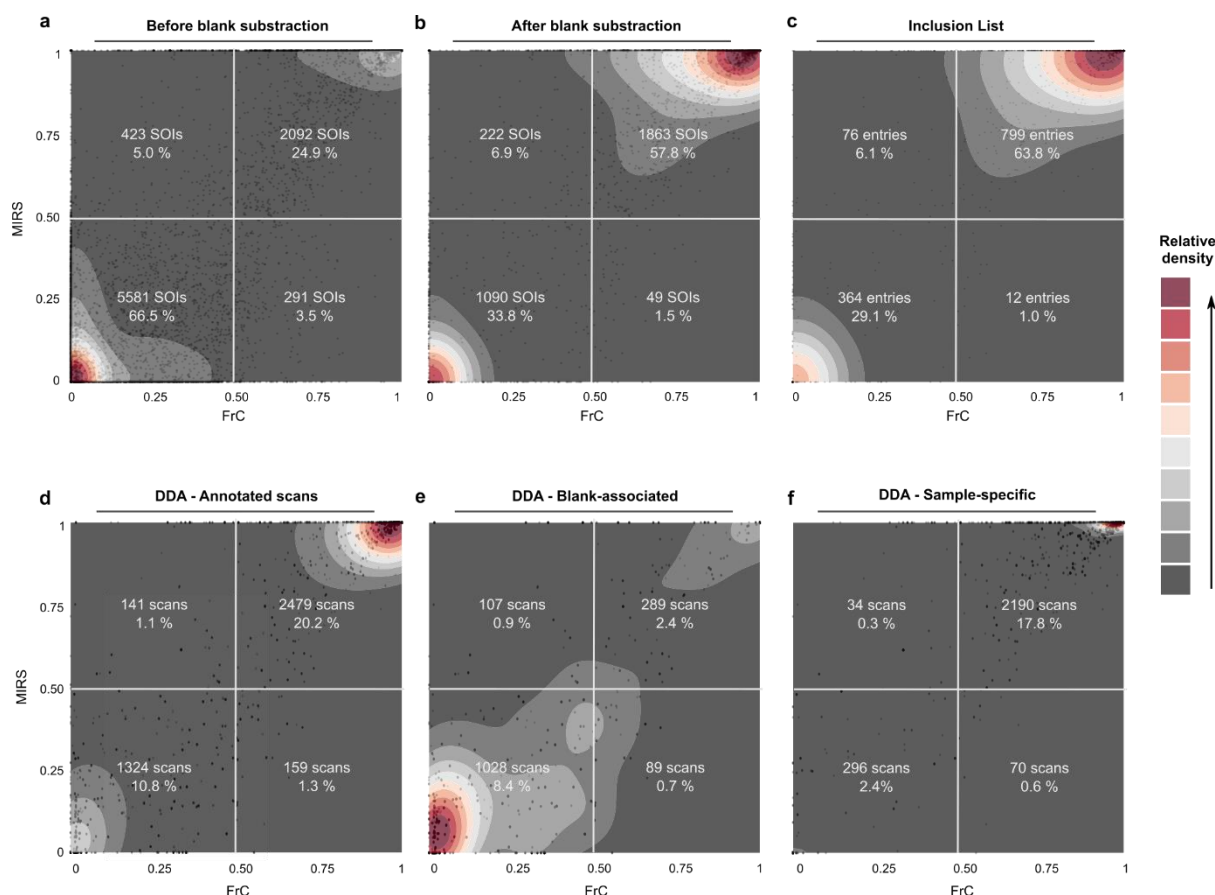
178
179

180 **Figure 3. Distribution of MS2 scans acquired by HERMES and iterative DDA.** a) Unlabelled *E. coli*
181 and b) human plasma samples acquired by iterative DDA. c) Human plasma sample acquired by
182 iterative DDA with background-exclusion. The acquired scans have been binned into 5Da-5s intervals.
183 The Venn diagrams show the bin intersections between HERMES and DDA. The precursor *m/z* of DDA
184 scans have been queried into the corresponding ionic formula *m/z* database with a 3 ppm mass error
185 tolerance. DDA scans annotated in the database were further classified according to whether the *m/z*
186 and retention time of the scans could be matched to the HERMES inclusion list or not. Percentages in
187 the pie-charts refer to the total number of acquired DDA MS2 scans.

188

189 To confirm the biogenic specificity of the MS2 scans in HERMES, a reference ¹³C-
190 labeled (at ≥98% from uniformly ¹³C-labeled glucose) *E. coli* credentialing extract was
191 analysed under identical LC/MS1 conditions. For each selected precursor ion in the unlabelled
192 *E. coli* sample, we calculated its fractional contribution (FrC)²¹⁻²³ and the monoisotopic ratio
193 score (MIRS) by using the analogue ¹³C-labeled sample (see Online Methods). A metabolite
194 with *n* carbon atoms can have zero (FrC=0) to *n* (FrC=1) of its carbon atoms labelled with ¹³C.
195 In turn, similar intensity of the monoisotopic ion in the unlabelled and ¹³C-labelled *E. coli*
196 extracts indicates no isotopic enrichment (MIRS=0), whereas loss of intensity in the ¹³C-
197 labelled sample is associated with enrichment (MIRS=1). Around 63% of inclusion list entries
198 in HERMES were associated with highly ¹³C-enriched metabolites (FrC and MIRS>0.5),

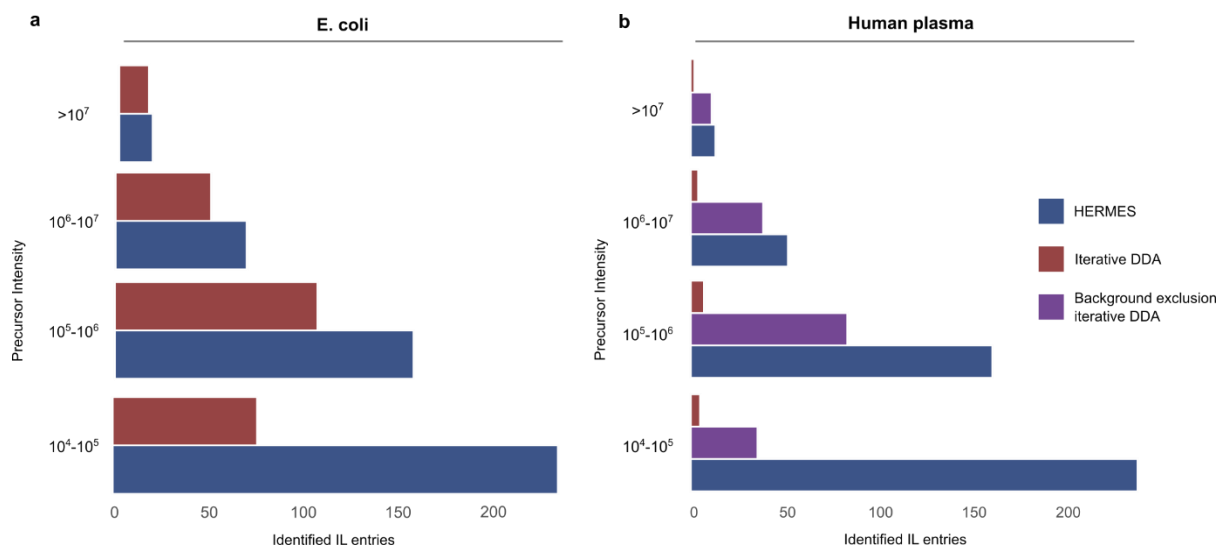
199 proving the biosynthetic origin of these ions (Fig. 4a-c)²⁴. These are mainly associated with
 200 abundant ions, while unlabelled precursors relate more frequently to low-abundant ions
 201 (Suppl. Fig. 11a,b). In contrast, only 20% of all DDA scans were associated with ¹³C-labeled
 202 and annotated precursors from ECMDB and the KEGG database, pointing to ions also present
 203 in the blank sample as the main source of unlabelled precursors (Fig. 4d-f). ¹³C-labeled
 204 precursors in DDA corresponded to highly abundant ions that were also covered by inclusion
 205 list entries in HERMES (Suppl. Fig. 11c).
 206



207
 208 **Figure 4. ¹³C-enrichment analysis in the labelled *E. coli* sample.** Each panel represents a scatterplot
 209 of two independent isotopic enrichment scores – FrC (Fractional Contribution) and MIRS (Monoisotopic
 210 Ratio Score) – and an overlaid density estimation. a) Distribution of SOIs before applying the blank
 211 subtraction filtering in HERMES. b) Same SOI list after removing most blank-related SOIs. c) SOIs in
 212 the MS2 inclusion list after removing redundant signals from b). d) Iterative DDA scans that could be
 213 matched to any *m/z* of the ionic formula database. e) DDA scans associated with SOIs removed during
 214 the blank subtraction step from a) to b). f) DDA scans associated with SOIs conserved during the blank
 215 subtraction. Percentages in a), b) and c) correspond to the total number of SOIs and inclusion list
 216 entries, accordingly, while percentages in d), e) and f) correspond to the total number of acquired DDA
 217 scans.

219 The biogenic specificity of HERMES resulted in higher similarity scores by mass
 220 spectral matching in databases (MassBankEU, MoNA, HMDB, Riken, NIST14, mzCloud)²⁵
 221 than iterative DDA (see Online Methods). HERMES provided nearly double the number of
 222 confident structural metabolite annotations than iterative DDA (Fig. 5a and Suppl. Fig. 12a).
 223 The higher identification rate of HERMES was validated by using alternative spectral similarity
 224 and distance metrics (Suppl. Fig. 13). A fraction of the ¹³C-labeled compounds, however, could
 225 not be identified due to low intensity SOIs and/or the lack of reference spectra in databases.
 226 For the former, setting the maximum ion injection time at high values (1,500 ms) improved
 227 sensitivity and MS2 spectral quality in HERMES, resulting in more informative fragments and
 228 better spectral matching (Suppl. Fig. 14). Furthermore, we identified unlabelled metabolites
 229 (FrC=0) in the ¹³C-labeled *E.coli* sample, such as choline, that we attribute to contaminants of
 230 the minimal growth medium that could not properly be removed by blank subtraction.

231



232

233 **Figure 5. Identified inclusion list entries according to the MS1 precursor intensity.** An inclusion
 234 list (IL) entry is considered identified if at least one MS2 scan associated with it has a compound hit in
 235 the reference MS2 database with either cosine score > 0.8 (in-house database from MassBankEU,
 236 MoNA, Riken and NIST14 spectra), or Match > 90 and Confidence > 30 (mzCloud). Positive ionization
 237 data. a) *E. coli* extract. b) Human plasma extract.

238

239 Finally, we used a human plasma extract to compare HERMES and iterative DDA, with
 240 and without background exclusion²⁶. Here we used 23,797 unique molecular formulas from
 241 the HMDB and Chemical Entities of Biological Interest (ChEBI) database to explore virtually
 242 all known exogenous and endogenous small molecules in this biofluid. HERMES generated
 243 110,387 and 46,973 ionic formulas that covered 60% and 14% of all data points acquired by

244 LC (RPC18)-Orbitrap MS in positive and negative ionization mode, respectively (Fig. 2b and
245 Suppl. Fig 9b). Consistent with the pattern observed for *E. coli*, 8% of all acquired data points
246 were associated with an XCMS peak (Fig. 2b), which resulted in less than 10% overlap of the
247 most frequent adducts annotations between HERMES and CliqueMS/CAMERA (Suppl. Fig.
248 10c,d). Only 1.2% of the data points in XCMS peaks matched with an ionic formula from HMDB
249 or ChEBI and were present in the inclusion list. Again, more than half of DDA precursors could
250 not be annotated as monoisotopic ionic formulas from HMDB and ChEBI without blank
251 subtraction (Fig. 3b). As expected, background exclusion in iterative DDA increased to 28%
252 the number of common MS2 scans between HERMES and DDA (Fig. 3c). Yet, the number of
253 confident structural metabolite identifications with HERMES was more than three times
254 greater than DDA because of the larger coverage of sample-specific and low abundant
255 precursor ions (Fig. 5b and Suppl. Fig. 12b).

256

257 **Discussion**

258 Our results demonstrate that a conventional LC/MS-based untargeted metabolomic
259 experiment can contain up to ~50 times more non-specific and redundant data points than
260 sample-specific and selective ones, which can account for as much as 90% of the MS2
261 acquisition run time in a state-of-the-art iterative DDA experiment. Current untargeted
262 metabolomic approaches are unable to properly annotate the large number of 'junk' MS1 and
263 MS2 signals, leading to low-quality MS2 spectra, false-positive identifications and an overall
264 low number of identified metabolites. HERMES solves this problem by implementing a broad
265 scope and molecular formula-oriented method that improves MS2 coverage by optimizing
266 MS2 acquisition time focusing on sample-specific, MS1 pre-annotated, and biologically
267 relevant compounds. In contrast to DDA, the continuous targeted MS2 acquisition in HERMES
268 does not acquire MS1 scans nor use dynamic exclusion, allowing a cleaner and easier
269 deconvolution and thereby increasing the quality of MS2 spectra and the number of identified
270 metabolites. This is also a differentiating element compared to DIA, where multiple precursor
271 ions, including redundant and biologically irrelevant ions, are simultaneously fragmented.

272 Although HERMES relies on high resolution MS data, it can be applied to the full range
273 of HRMS instrumentations and mass resolutions. We tested HERMES using an LC-QTOF MS
274 at ~30,000 resolution with the only drawback that certain isotopic patterns were difficult to
275 annotate due to lower isotopic fidelity and possible interferences, as shown in Suppl. Fig. 1.

276 The coverage of compounds detected by HERMES in any given sample is determined
277 by the list of molecular formulas. We currently provide ready-to-use lists for most important

278 databases in biomedical and/or environmental studies: HMDB, ChEBI, NORMAN, KEGG,
279 LipidMaps²⁷ and LipidBlast²⁸, which can be used individually or merged in a customized way,
280 leading to $>10^3$ - 10^4 unique molecular formulas. Yet, HERMES provides maximum
281 experimental flexibility by allowing users to add new molecular formulas not reported in public
282 databases²⁹, including *in silico* secondary metabolism prediction³⁰⁻³² such as environmental
283 microbial degradation, biotransformations of gut and soil/aquatic microbiota, or small peptides
284 such as dipeptides and tripeptides. Therefore, in a context where virtually all known molecular
285 formulas ($>10^3$ - 10^4) can be covered, we believe that ‘untargeted’ and ‘broad scope targeted’
286 are interchangeable terms to define HERMES.

287 Future developments should provide optimized maximum ion injection time and
288 collision energies for each inclusion list entry to reduce the number of MS2 scans required,
289 and improve the quality of MS2 spectra, particularly for low intensity SOIs. Also, the ability to
290 annotate in-source fragments by HERMES is determined by the size and content of public
291 (and in-house) MS2 databases. In-source fragments from metabolites that are not present in
292 a MS2 database, or do not contain low-energy MS2 data, will not be annotated. As a rough
293 guide, we currently filter out ~10% of the SOIs as in-source fragments, however, this
294 percentage may change as more metabolites are added to MS2 databases. Finally, the use
295 of sample-specific and high-quality MS2 spectra linked to preannotated precursor ions (i.e.,
296 molecular formula, adduct) should restrict the range of known and unknown chemical
297 structures for *in silico* MS2 fragmentation tools. This includes the possibility to find unknown
298 isomeric forms from known molecular formulas, provided that these novel structural isomers
299 produce sufficiently specific fragmentation spectra.

300

301 **Online Methods**

302 **Materials.** LC/MS-grade acetonitrile, water, isopropanol, and methanol (Burdick & Jackson)
303 were purchased from Honeywell (Muskegon, MI). LC/MS-grade ammonium bicarbonate,
304 ammonium hydroxide and methylenediphosphonic (medronic) acid were purchased from
305 Sigma-Aldrich (St. Louis, MO). Dried down metabolic extracts of *E. coli* were purchased from
306 Cambridge Isotope Laboratories (MSK-CRED-DD-KIT). Spike-in compounds (Suppl. Table 1)
307 were purchased from Sigma-Aldrich (Zwijndrecht, The Netherlands), LGC Standards (Wesen,
308 Germany) and Toronto Research Chemicals (Toronto, ON).

309

310 **Sample preparation**

311 **Environmental water.** Surface water was obtained from the Lekkanaal at Nieuwegein (The

312 Netherlands). The spike-in compounds were added to the surface water sample to a final
313 concentration of 1 µg/L. Subsequently, the sample was filtered using Phenex™ reversed
314 cellulose 15 mm Syringe Filters 0.2µ (Phenomenex, Torrance, USA) and transferred to a LC
315 autosampler vial.

316 **E.coli.** Dried down *E. coli extracts* (unlabelled and uniformly ¹³C-labelled) were reconstituted
317 in 100 µL of acetonitrile:water (2:1), followed by 30 s vortexing, 5 min of sonication, and 30 s
318 of vortexing.

319 **Human plasma.** Plasma aliquots (50 µL) were thawed at 4°C and briefly vortex-mixed.
320 Proteins were precipitated by the addition of 200 µL cold acetonitrile/methanol/water (5:4:1,
321 vol/vol) followed by 10 seconds vortex-mixing. Samples were subsequently maintained on ice
322 for 30 min. After centrifugation (10 min, 15.200 rpm at 4°C), 100 µL of supernatant were
323 transferred to a LC autosampler vial.

324

325 **LC-MS analysis**

326 **Environmental water and human plasma.** Ultra-high performance LC (UHPLC)/MS was
327 performed with a Thermo Scientific Vanquish UHPLC system interfaced with a Thermo
328 Scientific Orbitrap Fusion Tribrid mass spectrometer operated in positive or negative ion
329 mode. Reverse phase C18 liquid chromatography (RPLC) analysis was performed by using a
330 Xbridge BEH C18 column (Waters, Etten-Leur, The Netherlands) with the following
331 specifications: 150 mm x 2.1 mm, 2.5 µm. Mobile-phase solvents were composed of A =
332 ultrapure water with 0.05% formic acid (v/v) and B = acetonitrile with 0.05% formic acid (v/v).
333 The column compartment was maintained at 25 °C for all experiments. The following linear
334 gradient was applied at a flow rate of 250 µL/min: 0-1 min: 5% B, 1-25 min: 5-100% B, 25-29
335 min: 100% B, 29.0-29.5 min 5% B followed by 4.5 min of re-equilibration phase. One µL of the
336 human plasma extract was diluted in 100 µL of ultrapure water, and the injection volume was
337 100 µL for all experiments. Data were collected with the following settings: spray voltage, 3.0
338 kV and -2.5 kV in positive and negative mode, respectively; sheath gas, 40; auxiliary gas, 10;
339 sweep gas, 5; ion transfer tube temperature, 300 °C; vaporizer temperature, 300 °C; mass
340 range, 80-1000 Da; RF lens, 50%; resolution, 120,000 (MS1), 15,000 (MS2); AGC target, 2e5
341 (MS1), 5e4 (MS2); maximum injection time, 100 ms (MS1), 50 ms (HERMES), 50 ms (DDA);
342 isolation window, 1.6 Da. The collision energy was 35% for HCD fragmentation. With every
343 batch run, mass calibration was performed using Pierce ESI positive and negative ion
344 calibration solution to obtain a mass error of <2 ppm.

345 **E.coli.** LC/MS was performed with a Thermo Scientific Vanquish Horizon UHPLC system
346 interfaced with a Thermo Scientific Orbitrap ID-X Tribrid Mass Spectrometer (Waltham, MA).
347 Hydrophilic interaction liquid chromatography (HILIC) analysis was performed by using a

348 SeQuant ZIC-pHILIC column (Merck Millipore, Burlington, MA) with the following
349 specifications: 150 mm x 2.1 mm, 5 μ m. Mobile-phase solvents were composed of A = 20 mM
350 ammonium bicarbonate, 0.1% ammonium hydroxide solution (25% ammonia in water) and 2.5
351 μ M medronic acid in water:acetonitrile (95:5) and B = 95% acetonitrile, 5% water, 2.5 μ M
352 medronic acid. The column compartment was maintained at 40 °C for all experiments. The
353 following linear gradient was applied at a flow rate of 250 μ L min⁻¹: 0-1 min: 90% B, 1-12 min:
354 90-35% B, 12.5-14.5 min: 25% B, 15 min: 90% B followed by 4 min of re-equilibration phase
355 at 400 μ L min⁻¹ and 2 min at 250 μ L min⁻¹. The injection volume was 2 μ L for all experiments.
356 Data were collected with the following settings: spray voltage, 3.5 kV and -2.8 kV in positive
357 and negative mode, respectively; sheath gas, 50; auxiliary gas, 10; sweep gas, 1; ion transfer
358 tube temperature, 300 °C; vaporizer temperature, 200 °C; mass range, 70-1000 Da; RF lens,
359 60%; resolution, 120,000 (MS1), 15,000 (MS2); AGC target, 2e5 (MS1), 5e4 (MS2); maximum
360 injection time, 200 ms (MS1), 35 ms (HERMES, unless otherwise stated), 100 ms (iterative
361 DDA); isolation window, 1 Da. The collision energy was 35% for HCD fragmentation.

362 **Iterative DDA**

363 **E.coli.** After the first DDA run, the raw data file containing MS2 spectra was converted to an
364 .MS2 file using MS Convert³³ Next, the IEomics tool³⁴ was used to generate the first exclusion
365 list of features fragmented in the first DDA run. User inputs in the R script were RTWindow =
366 0.3 min, noiseCount = 25, MZWindow = 0.001. This procedure was repeated two times, which
367 resulted in a total of three DDA data runs per polarity. The mass tolerance for exclusion lists
368 was 5 ppm.

369 **Plasma.** An exclusion list of background ions was generated using the AcquireX workflow of
370 Xcalibur data acquisition software (Thermo Fisher Scientific), by analysing an ultrapure water
371 sample. The exclusion list contains the exact mass, retention window and intensity (exclusion
372 override factor = 3) of the excluded background ions. DDA was performed for the top 6-8
373 most intense ions per full scan. Dynamic exclusion was used to prevent redundant acquisition
374 of MS2 spectra for a selected precursor ion for 10 s, when two MS2 spectra were acquired
375 within 20 s, resulting in a total of three DDA data runs per polarity. A mass tolerance of 5 ppm
376 was used for the exclusion list and dynamic exclusion.

377

378 **HERMES algorithm**

379 All analysis were performed using RHermes (version 0.99.0).

380 **MS1 data processing:** Theoretical isotopic patterns of each ionic formula were calculated by
381 Envipat (version 2.4) and refined by RHermes, based on the predefined experimental mass

382 resolution and mass accuracy values. Local resolution was calculated for each ionic formula

383 as: $R(mz) = R_{ref} \cdot \sqrt{\frac{mz}{mz_{ref}}}$.

384 Using as input a set of mzML files, SOIs were detected by RHermes using two sets of 5s bins
385 (offset by 2.5s) and required a minimum scan density of 30% of acquired scans.

386 Blank subtraction was performed using an heuristic prefilter (intensity ratio sample/blank > 3)
387 and an artificial neural network (ANN) trained with >3000 manually annotated sample/blank
388 SOI comparisons from two different biological matrices. Each blank-sample pair was manually
389 classified according to whether the elution shape and intensity were similar (0) or not (1). The
390 entries were then separated into 80/20 training and testing sets. The model was trained using
391 Keras (<https://keras.io/>), with the categorical crossentropy metric as a performance indicator.
392 After XXXXX rounds of training, the model obtained a XX% accuracy in the testing set. Adduct
393 and isotopologue grouping were performed using a cosine shape similarity score and required
394 a cosine >0.8. Cosine similarity has a superior discriminatory power than the Pearson
395 correlation¹³ for discriminating between features corresponding to the same metabolite from
396 coeluting features corresponding to different metabolites¹⁵.

397 In-source fragment (ISF) annotation was performed using an in-house MS2 database
398 consisting of 675,663 low-energy spectra (<20% HCD, <10eV CID) from MassBankEU,
399 MoNA, HMDB, Riken and NIST14. Low-energy spectra were selected according to each SOI
400 formula annotation. Intense (>20% of maximum intensity) fragments m/z were then queried
401 against the SOI list. Finally, the suspected ISF SOIs elution profiles were compared to the
402 original SOI and a cosine similarity score was calculated.

403 **MS2 data processing.** The program exports the inclusion list into a csv file used to generate
404 the MS2 acquisition method. Acquired MS2 scans were linked to each inclusion list entry; if
405 >5 scans were acquired, a deconvolution algorithm was applied, where fragments m/z were
406 grouped and initially split with a Centwave peak picking (peakwidth = c(5,60)). A cosine shape
407 similarity score was applied to each pair of fragment peaks to generate a similarity network.
408 Each network was then partitioned using a greedy algorithm from *igraph* (version 1.2.4.2) and
409 yielded a variable number of deconvoluted MS2 spectra (see Algorithm 3). If fewer than 5
410 scans were acquired, the scan with the highest TIC was selected and filtered by intensity (>
411 0.5% of maximum).

412 Note: the program can also process conventional DDA data to filter out blank-related MS2
413 scans. RHermes uses DDA's MS1 scans to generate a SOI list and uses only those MS2
414 scans that match with the SOI list.

415

416

417 **XCMS data processing**

418 LC/MS1 raw data files (ESI+ and ESI- modes) were converted to open standard format mzML
419 using Proteowizard MS-convert³³ and subsequently processed by HERMES and XCMS
420 software¹⁹ (version 3.8.1). XCMS settings were: `xcmsSet(method="centWave", ppm=3,`
421 `peakwidth=c(10, 60))`; Common data points between SOIs in HERMES and XCMS peaks were
422 calculated by extracting the raw data points delimited by each XCMS peak ($rt_{\min} < rt < rt_{\max}$
423 and $mz_{\min} < mz < mz_{\max}$) and generating the set intersections using `dplyr` (version 1.0.4).

425 **CAMERA and CliqueMS data processing**

426 Both CAMERA and CliqueMS were run using their respective default adduct lists, which
427 contained all adducts considered by HERMES. Mass error for both tools was set to 3 ppm.

429 **Uniformly ¹³C-labeled *E. coli***

430 Fractional contribution (FrC) was calculated using the formula:

$$431 \quad FrC = \sum_{i=1}^N \frac{M_i \cdot i}{M_{O_{unlab}} \cdot n}$$

432 where N is the number of carbon atoms in the molecule.

433 Monoisotopic Ratio Score (MIRS) was calculated using the formula:

$$434 \quad MIRS = 1 - \frac{M_{O_{labelled}}}{M_{O_{unlabelled}}}$$

435 If MIRS was smaller than zero, it was set to zero so that all points range from 0 to 1.

437 **MS2 identification**

438 **In-house DB.** MS2 spectra were obtained from MassBankEU, MoNA, HMDB, Riken and
439 NIST14 databases. RHermes includes the MassBank EU database for testing purposes,
440 which is available as an .rds file at Zenodo (accession number 4678268).

441 All fragment m/z were discretized into 0.01Da bins. Each spectrum precursor m/z was
442 matched against the DB spectra m/z with a 0.01Da tolerance. For the HERMES matching, the
443 reference spectra were further filtered according to the formula database used in the MS1
444 analysis. A cosine similarity score was calculated between the query and reference spectra
445 and resulting hits were filtered by requiring a score > 0.8 .

446 **mzCloud DB.** The processed HERMES MS2 spectra were exported to the mzML file format.
447 The DDA files were directly imported through MassFrontier version 8.0 SR1 (Thermo
448 Scientific) and matched against the mzCloud database using three component identification
449 types: Identity, Similarity Forward and Similarity Reverse; with the following constraints: 4.0

450 Tolerance Factor and Match Ion Activation Type. The resulting hits were filtered by both Match
451 and Confidence scores (requiring a score > 90 and > 30, respectively).

452 Identified inclusion list entries (Figure 5 and Supp Figure 11) were calculated as number of
453 inclusion list entries that resulted in a valid hit (i.e. high score) against either of the two
454 databases. For DDA, this number was calculated by matching the precursor *m/z* and RT of
455 the scans to the inclusion list and then examining if (i) any of the scans have at least one valid
456 hit against either of the two databases and (ii) any valid hit had a molecular formula present
457 in the HERMES formula database.

458 All similarity metrics were calculated using the R package *philterropy* (version 0.4.0). MS2
459 spectra were discretized into 0.01Da bins and their fragment intensities scaled by the sum of
460 the intensities, so that all calculated metrics were comparable across the spectra. The query
461 spectra (both DDA and HERMES) were matched against the previously described In-house
462 DB. For each query, all DB hits were grouped, taking the maximum similarity (cosine and
463 fidelity) and the lowest distance (squared chord and topsoe). Additionally, HERMES hits were
464 restricted to compounds with formulas present in the HERMES formula database. The
465 corresponding plots were generated using *ggplot2* (version 3.3.3).

466

467 **Data availability**

468 Input mzML/mzXML mass spectrometry data files and an RMarkdown are available at Zenodo
469 with the accession number 4678268.

470

471 **Code availability**

472 The source code of RHermes is offered to the public as a freely accessible software package
473 under the GNU GPL, version 3 license, and is available at
474 <https://github.com/RogerGinBer/RHermes>.

475

476 **Acknowledgements**

477 We gratefully acknowledge financial support by Ministerio de Educación y Formación
478 Profesional (Spanish Government) to R.G. (2020-COLAB-00552). O.Y. was supported by
479 Ministerio de Economía y Competitividad (MINECO) (BFU2014-57466-P), Spanish
480 Biomedical Research Centre in Diabetes and Associated Metabolic Disorders (CIBERDEM),
481 an initiative of Instituto de Investigación Carlos III (ISCIII), and the European Union's Horizon
482 2020 program (MSCA-ITN-2015; 675610). We thank members of the Mil@b for helpful
483 comments.

484

485

486 **Author contributions**

487 RG and OY designed the research. RG, JC, JMB, MV and OY developed the computational
488 method. DV and MSH performed LC-MS and MS2 experiments. All authors applied and
489 evaluated the method on biological samples. RG and OY wrote the manuscript, in cooperation
490 with all authors.

491

492 **Competing interests**

493 The authors declare no personal financial interests. A patent application for the method has
494 been filled by RG, JC and OY (P202030061).

495

496

497 **References**

498 1. Sindelar, M. & Patti, G. J. Chemical Discovery in the Era of Metabolomics. *J. Am. Chem.*
499 *Soc.* **142**, 9097–9105 (2020).

500 2. Duan, L., Molnár, I., Snyder, J. H., Shen, G. & Qi, X. Discrimination and Quantification of
501 True Biological Signals in Metabolomics Analysis Based on Liquid Chromatography-
502 Mass Spectrometry. *Mol. Plant* **9**, 1217–1220 (2016).

503 3. Myers, O. D., Sumner, S. J., Li, S., Barnes, S. & Du, X. Detailed Investigation and
504 Comparison of the XCMS and MZmine 2 Chromatogram Construction and
505 Chromatographic Peak Detection Methods for Preprocessing Mass Spectrometry
506 Metabolomics Data. *Anal. Chem.* **89**, 8689–8695 (2017).

507 4. Domingo-Almenara, X., Montenegro-Burke, J. R., Benton, H. P. & Siuzdak, G.
508 Annotation: A Computational Solution for Streamlining Metabolomics Analysis. *Anal.*
509 *Chem.* **90**, 480–489 (2018).

510 5. Tsugawa, H. *et al.* MS-DIAL: data-independent MS/MS deconvolution for comprehensive
511 metabolome analysis. *Nat. Methods* **12**, 523–526 (2015).

512 6. Yin, Y., Wang, R., Cai, Y., Wang, Z. & Zhu, Z.-J. DecoMetDIA: Deconvolution of
513 Multiplexed MS/MS Spectra for Metabolite Identification in SWATH-MS-Based
514 Untargeted Metabolomics. *Anal. Chem.* **91**, 11897–11904 (2019).

- 515 7. Guo, J. & Huan, T. Comparison of Full-Scan, Data-Dependent, and Data-Independent
516 Acquisition Modes in Liquid Chromatography–Mass Spectrometry Based Untargeted
517 Metabolomics. *Anal. Chem.* **92**, 8072–8080 (2020).
- 518 8. Röst, H. L. *et al.* OpenMS: a flexible open-source software platform for mass
519 spectrometry data analysis. *Nat. Methods* **13**, 741–748 (2016).
- 520 9. Huan, T. *et al.* Systems biology guided by XCMS Online metabolomics. *Nat. Methods*
521 **14**, 461–462 (2017).
- 522 10. Wishart, D. S. *et al.* HMDB 4.0: the human metabolome database for 2018. *Nucleic*
523 *Acids Res.* **46**, D608–D617 (2018).
- 524 11. J, H. *et al.* ChEBI in 2016: Improved services and an expanding collection of
525 metabolites. *Nucleic Acids Res.* **44**, D1214-9 (2015).
- 526 12. NORMAN Network *et al.* S0 | SUSDAT | Merged NORMAN Suspect List: SusDat. (2020)
527 doi:10.5281/zenodo.4249026.
- 528 13. Palmer, A. *et al.* FDR-controlled metabolite annotation for high-resolution imaging mass
529 spectrometry. *Nat. Methods* **14**, 57–60 (2017).
- 530 14. Domingo-Almenara, X. *et al.* Autonomous METLIN-Guided In-source Fragment
531 Annotation for Untargeted Metabolomics. *Anal. Chem.* **91**, 3246–3253 (2019).
- 532 15. Senan, O. *et al.* CliqueMS: a computational tool for annotating in-source metabolite ions
533 from LC-MS untargeted metabolomics data based on a coelution similarity network.
534 *Bioinformatics* **35**, 4089–4097 (2019).
- 535 16. Dührkop, K. *et al.* SIRIUS 4: a rapid tool for turning tandem mass spectra into metabolite
536 structure information. *Nat. Methods* **16**, 299–302 (2019).
- 537 17. Dührkop, K. *et al.* Systematic classification of unknown metabolites using high-resolution
538 fragmentation mass spectra. *Nat. Biotechnol.* 1–10 (2020) doi:10.1038/s41587-020-
539 0740-8.
- 540 18. Aron, A. T. *et al.* Reproducible molecular networking of untargeted mass spectrometry
541 data using GNPS. *Nat. Protoc.* **15**, 1954–1991 (2020).

- 542 19. Smith, C. A., Want, E. J., O'Maille, G., Abagyan, R. & Siuzdak, G. XCMS: Processing
543 Mass Spectrometry Data for Metabolite Profiling Using Nonlinear Peak Alignment,
544 Matching, and Identification. *Anal. Chem.* **78**, 779–787 (2006).
- 545 20. Kuhl, C., Tautenhahn, R., Böttcher, C., Larson, T. R. & Neumann, S. CAMERA: An
546 Integrated Strategy for Compound Spectra Extraction and Annotation of Liquid
547 Chromatography/Mass Spectrometry Data Sets. *Anal. Chem.* **84**, 283–289 (2012).
- 548 21. Buescher, J. M. *et al.* A roadmap for interpreting ¹³C metabolite labeling patterns from
549 cells. *Curr. Opin. Biotechnol.* **34**, 189–201 (2015).
- 550 22. Zamboni, N., Saghatelian, A. & Patti, G. J. Defining the Metabolome: Size, Flux, and
551 Regulation. *Mol. Cell* **58**, 699–706 (2015).
- 552 23. Jang, C., Chen, L. & Rabinowitz, J. D. Metabolomics and Isotope Tracing. *Cell* **173**,
553 822–837 (2018).
- 554 24. Mahieu, N. G., Huang, X., Chen, Y.-J. & Patti, G. J. Credentialing Features: A Platform
555 to Benchmark and Optimize Untargeted Metabolomic Methods. *Anal. Chem.* **86**, 9583–
556 9589 (2014).
- 557 25. Vinaixa, M. *et al.* Mass spectral databases for LC/MS- and GC/MS-based metabolomics:
558 State of the field and future prospects. *TrAC Trends Anal. Chem.* **78**, 23–35 (2016).
- 559 26. Cho, K. *et al.* Targeting unique biological signals on the fly to improve MS/MS coverage
560 and identification efficiency in metabolomics. *Anal. Chim. Acta* **1149**, 338210 (2021).
- 561 27. Fahy, E., Sud, M., Cotter, D. & Subramaniam, S. LIPID MAPS online tools for lipid
562 research. *Nucleic Acids Res.* **35**, W606-612 (2007).
- 563 28. Kind, T. *et al.* LipidBlast in silico tandem mass spectrometry database for lipid
564 identification. *Nat. Methods* **10**, 755–758 (2013).
- 565 29. Ludwig, M. *et al.* Database-independent molecular formula annotation using Gibbs
566 sampling through ZODIAC. *Nat. Mach. Intell.* **2**, 629–641 (2020).

567 30. Djoumbou-Feunang, Y. *et al.* BioTransformer: a comprehensive computational tool for
568 small molecule metabolism prediction and metabolite identification. *J. Cheminformatics*
569 **11**, 2 (2019).

570 31. Rutz, A. *et al.* Open Natural Products Research: Curation and Dissemination of
571 Biological Occurrences of Chemical Structures through Wikidata. *bioRxiv*
572 2021.02.28.433265 (2021) doi:10.1101/2021.02.28.433265.

573 32. Blin, K. *et al.* antiSMASH 5.0: updates to the secondary metabolite genome mining
574 pipeline. *Nucleic Acids Res.* **47**, W81–W87 (2019).

575 33. Chambers, M. C. *et al.* A cross-platform toolkit for mass spectrometry and proteomics.
576 *Nat. Biotechnol.* **30**, 918–920 (2012).

577 34. Koelmel, J. P. *et al.* Expanding Lipidome Coverage Using LC-MS/MS Data-Dependent
578 Acquisition with Automated Exclusion List Generation. *J. Am. Soc. Mass Spectrom.* **28**,
579 908–917 (2017).

580

581

582

583

584

585

586

587

588

589

590

591

592

Supplementary File

593
594

595

596 **HERMES: a molecular formula-oriented method to target the metabolome**

597

598

599

600 Roger Giné¹, Jordi Capellades^{1,2}, Josep M. Badia^{1,2}, Dennis Vughs³, Michaela Schwaiger-
601 Haber^{4,5}, Theodore Alexandrov^{6,7,8}, Maria Vinaixa^{1,2}, Andrea M. Brunner³, Gary J. Patti^{4,5},
602 Oscar Yanes*^{1,2}

603

604 1. Universitat Rovira i Virgili, Department of Electronic Engineering & IISPV, Tarragona,
605 Spain.

606 2. CIBER de Diabetes y Enfermedades Metabólicas Asociadas (CIBERDEM), Instituto
607 de Salud Carlos III, Madrid, Spain.

608 3. KWR Water Research Institute, Nieuwegein, The Netherlands.

609 4. Department of Chemistry, Washington University, St. Louis, MO, USA.

610 5. Department of Medicine, Washington University, St. Louis, MO, USA.

611 6. Structural and Computational Biology Unit, European Molecular Biology Laboratory,
612 Heidelberg, Germany.

613 7. Molecular Medicine Partnership Unit, European Molecular Biology Laboratory,
614 Heidelberg, Germany.

615 8. Skaggs School of Pharmacy and Pharmaceutical Sciences, University of California
616 San Diego, La Jolla, California, USA.

617

618

619

620 *Corresponding author:

621 Oscar Yanes, PhD

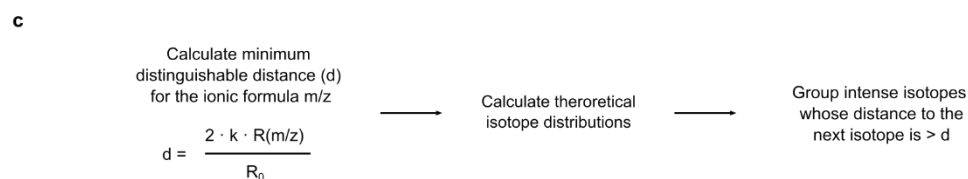
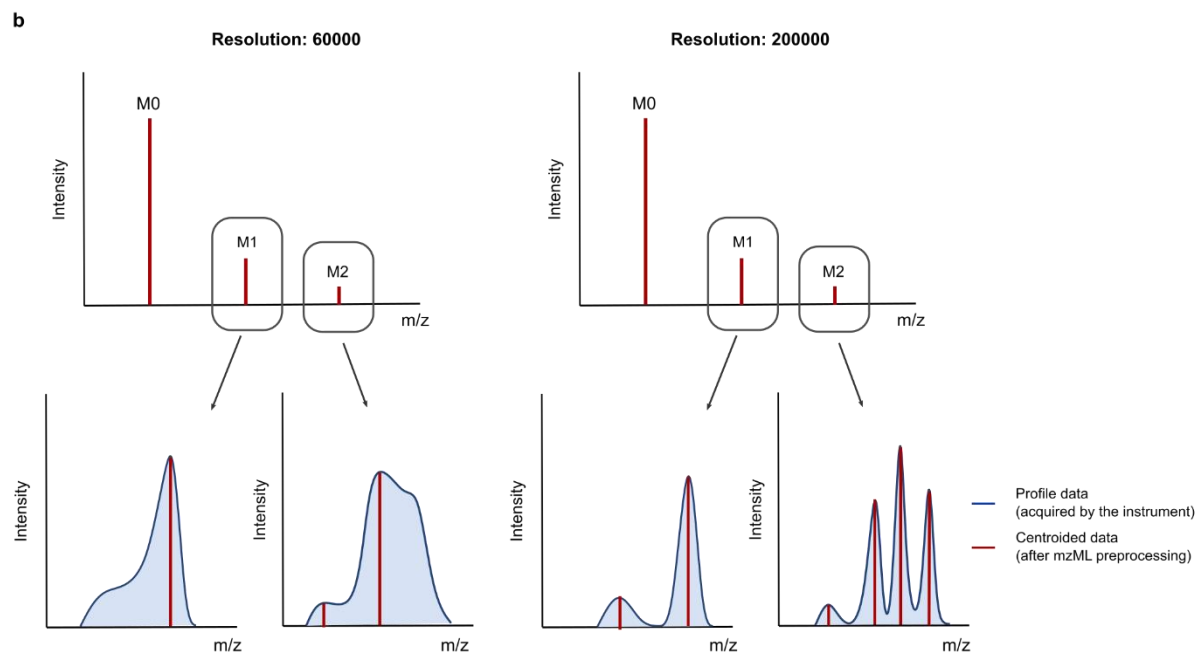
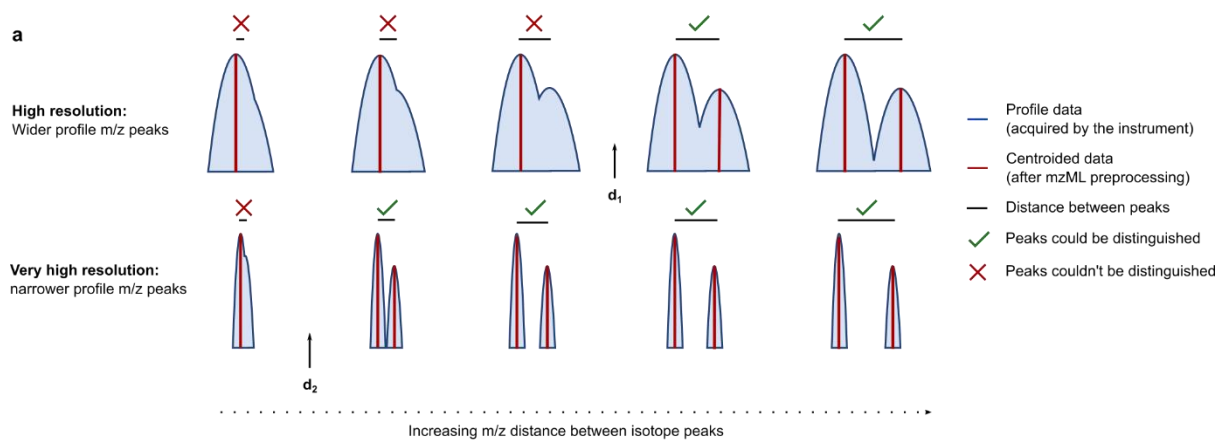
622 Department of Electronic Engineering

623 Universitat Rovira i Virgili

624 Avinguda Països Catalans, 26, 43007 Tarragona, Spain

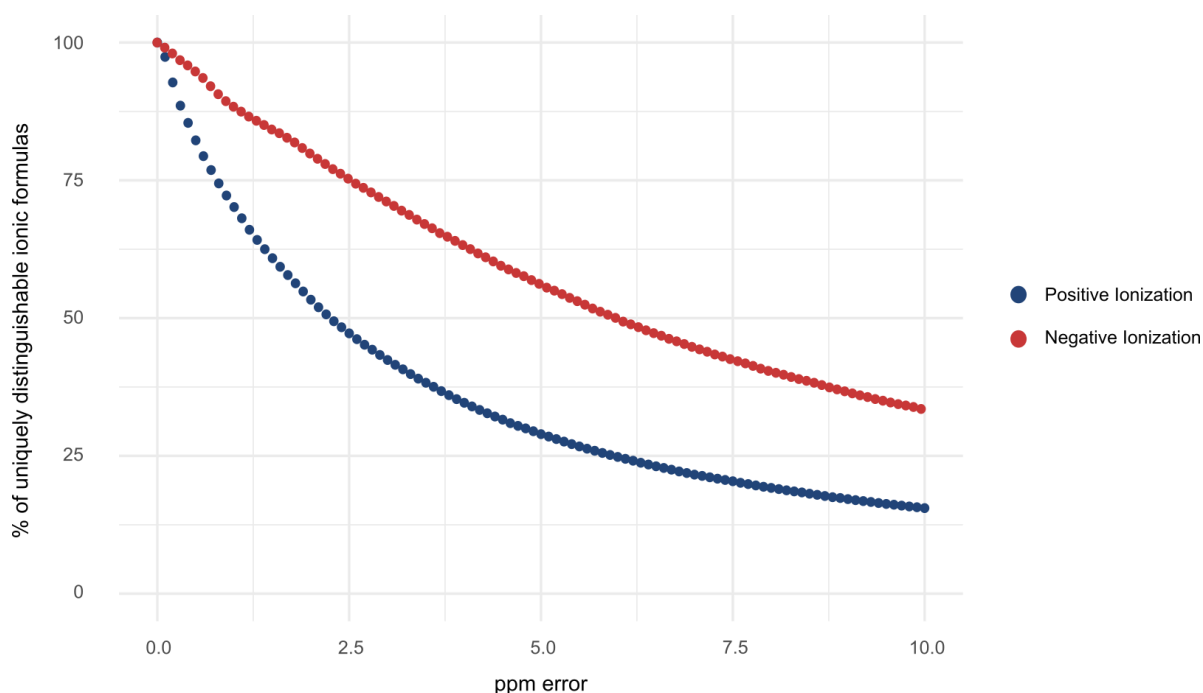
625 phone: +34 977759397

626 email: oscar.yanes@urv.cat

628 **Supplementary Figures**

629

630 **Supplementary Figure S1. Resolution-based isotopic envelope calculation.** a) The MS
 631 resolution is inversely proportional to the peak width of the acquired signals. When
 632 preprocessing raw MS1 data, a centroidization algorithm performs a peak-picking on a
 633 continuous profile signal (blue), yielding discrete, centroided signals (red). As resolution
 634 increases, the minimal distance to distinguish two adjacent peaks decreases ($d_2 < d_1$). b) In
 635 practice, this implies that, when acquiring data at lower resolutions, certain isotope signals are
 636 masked by close, more intense signals c) By calculating a resolution-based parameter d ,
 637 HERMES can estimate which close isotopologues can be distinguished in the acquired profile
 638 MS1 data and therefore present in the centroided data (See Algorithm 1).

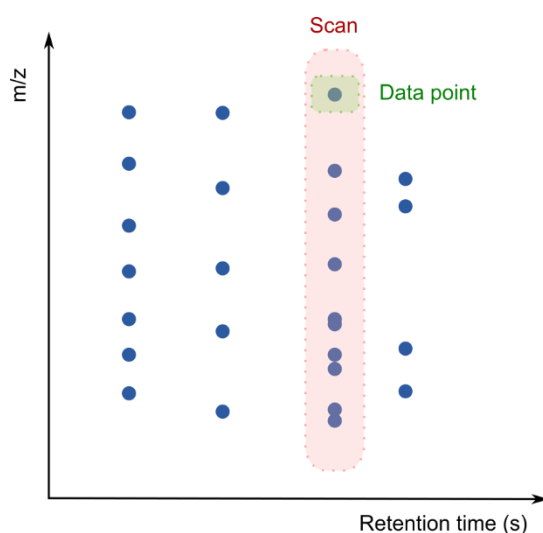


640

641 **Supplementary Figure S2. Distribution of uniquely distinguishable ionic formulas:** Ionic
 642 formula collisions from the NORMAN database (24,696 unique molecular formulas). Blue:
 643 Positive ionization taking $[M+H]^+$, $[M+Na]^+$, $[M+K]^+$, $[M+NH_4]^+$ and $[M]^+$ adducts. Red: Negative
 644 ionization taking $[M-H]^-$ and $[M+Cl]^-$ adducts. As the ppm error of the instrument increases, the
 645 larger the percentage of overlapping ionic formulas. The shape of the curve is closely related
 646 to the number of possible ionic formulas (distinct combinations of formulas-adducts) that are
 647 considered: the more ionic formulas, the larger the probability of overlaps.

648

649

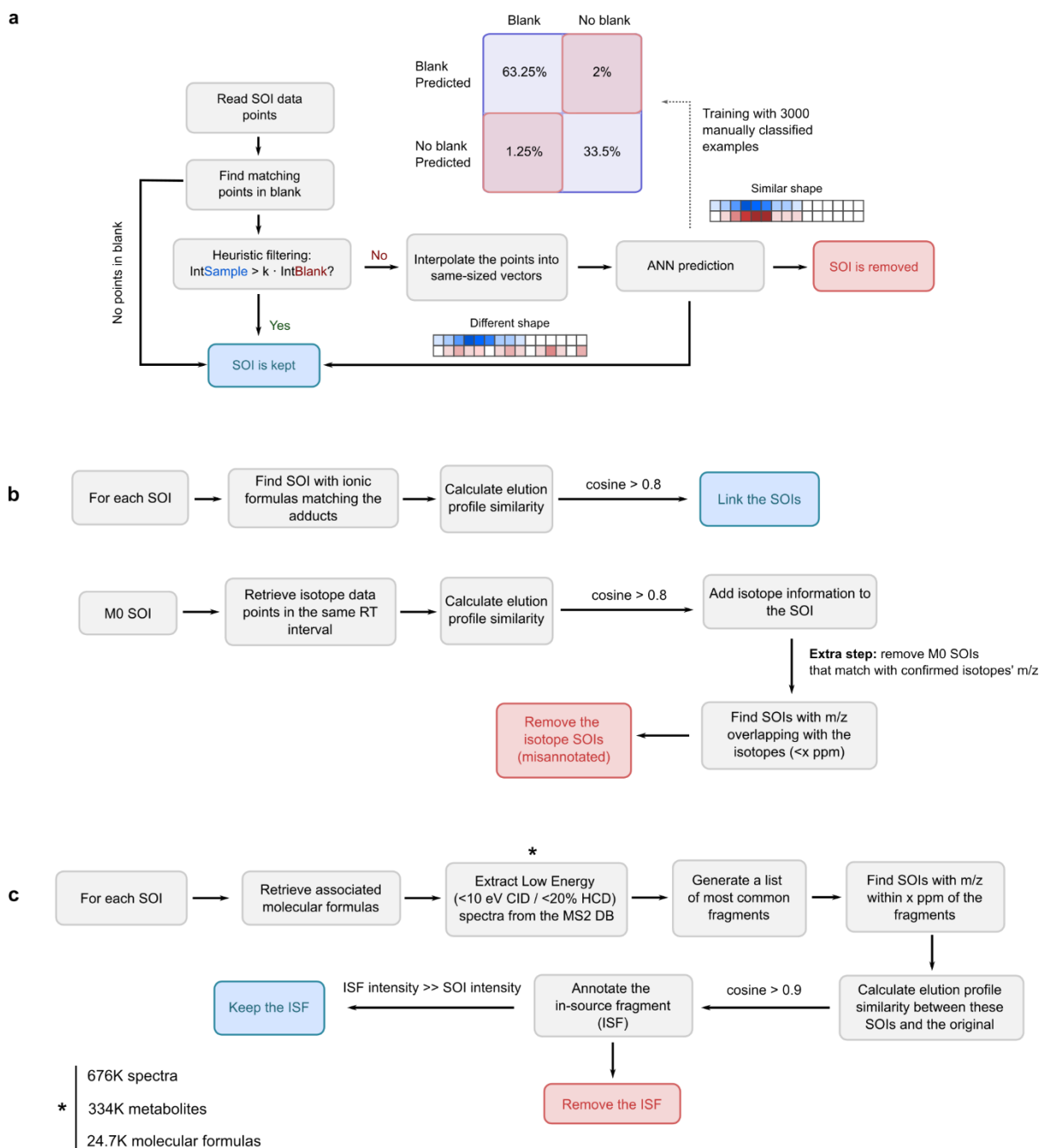


650

651 **Supplementary Figure S3: Difference between the terms scan and data point.** A *data*
 652 *point* is defined as a triplet $\{m/z, \text{Retention Time}, \text{Intensity}\}$, while a *scan* is the set of data
 653 points with the same retention time. A different concept is that of a *feature*, which is an abstract

654 concept defined as the m/z and retention time interval (composed of multiple scans) where a
 655 peak has been detected in one or multiple samples.

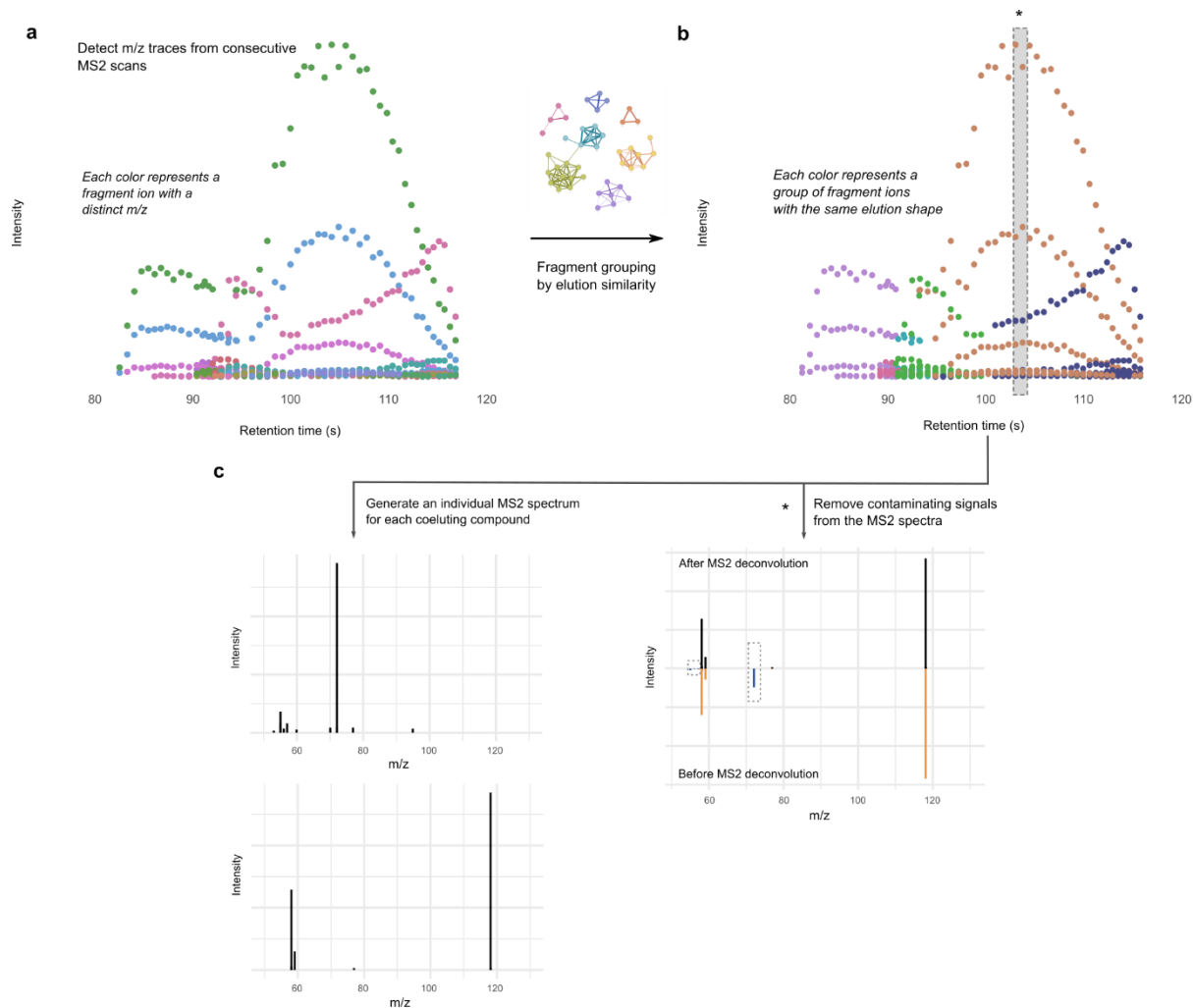
656



657

658 **Supplementary Figure S4. Schematic workflow of the different filtering steps in Hermes.**

659 a) Artificial neural network (ANN) for blank subtraction. b) Adduct and isotopologue grouping
 660 according to the similarity of their elution profiles. c) In-source fragment annotation, by using
 661 publicly available low-energy MS/MS data.

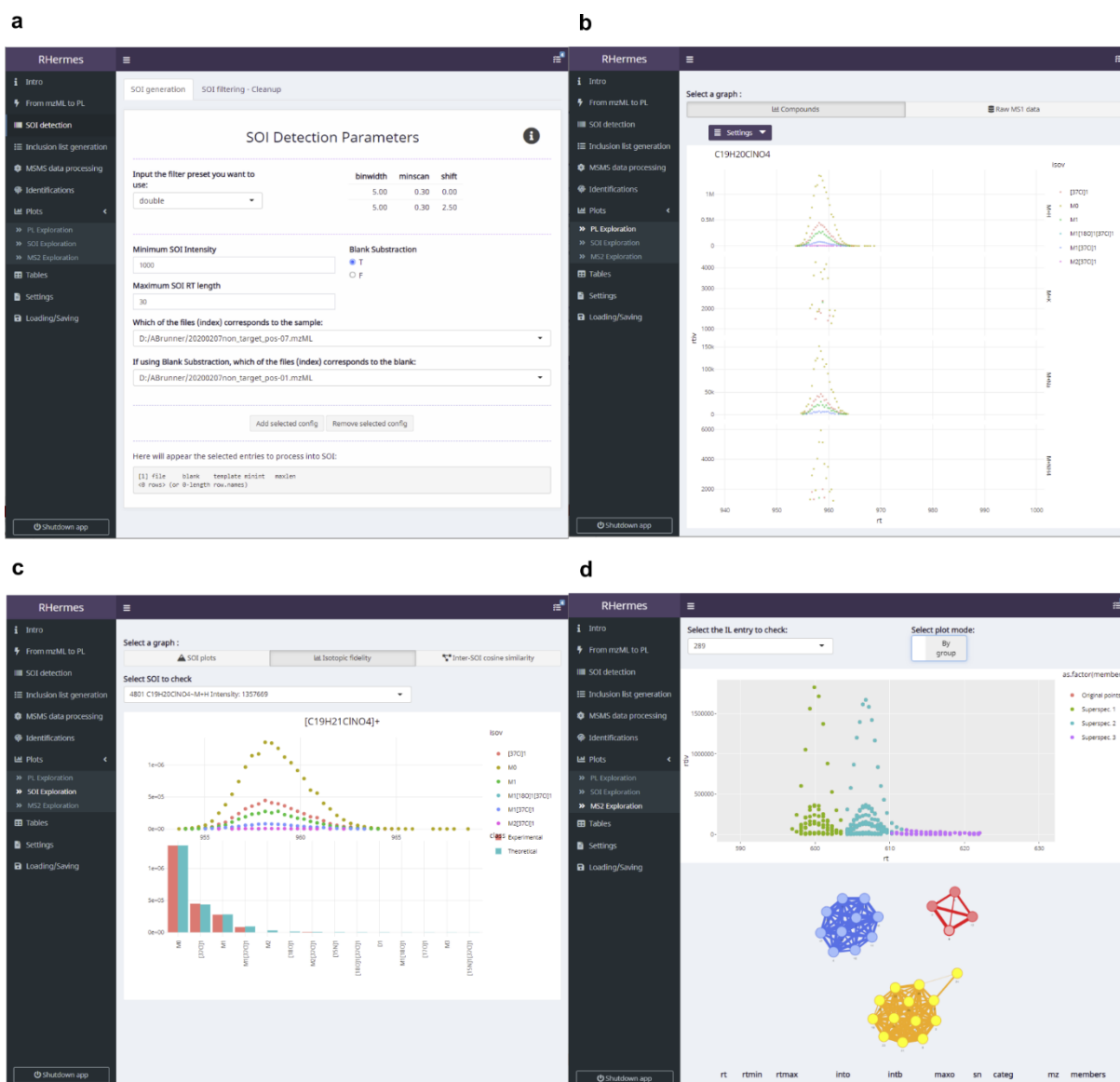


662

663

664 **Supplementary Figure S5. Continuous MS2 acquisition resolves co-eluting ionic**
 665 **species by comparing their fragment elution profile.** a) All fragment ions from continuous
 666 MS2 scans are grouped according to their m/z. b) A loose peak-picking algorithm is applied
 667 and the resulting peaks are grouped according to their elution profiles, generating a similarity
 668 network that is split by a greedy clustering algorithm. c) This grouping yields a curated MS2
 669 spectra for each coeluting species (see Algorithm 3). (*) The shaded slice shows the impact
 670 of the algorithm on the resulting MS2 spectral quality. The delineated fragments in blue have
 671 a different elution pattern from the rest and would contaminate the MS2 spectra if only one
 672 scan was acquired at the top of the peak. The grouping performed by HERMES confidently
 673 removes the contaminant ions and separates each group of fragments according to their
 674 elution.

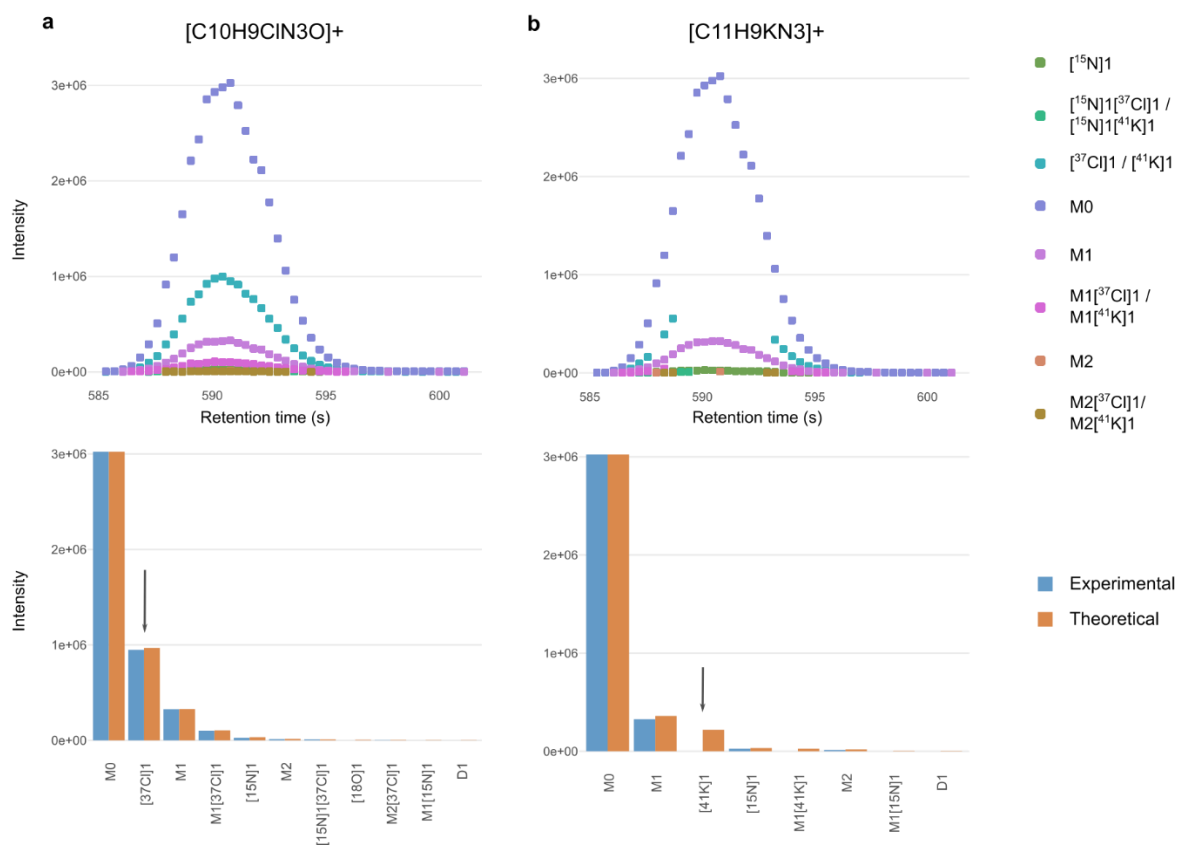
675



676

677

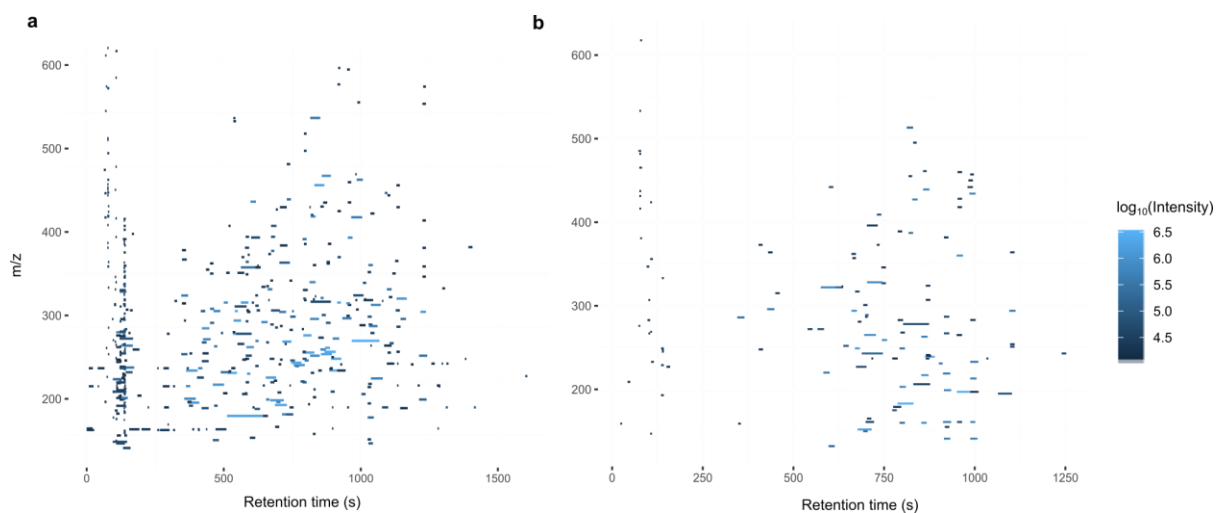
678 **Supplementary Figure S6. HERMES R Graphical User Interface (GUI).** a) Point-and-click
 679 selection of SOI detection parameters, with detailed explanations on their usage and optimal
 680 values. b) Visualization of isotopic profiles of different adducts of the same formula. Formulas
 681 can be inputted directly or inferred from the name of a selected compound. c) Isotopic fidelity
 682 exploration of selected SOIs. d) Visualization of the continuous MS2 deconvolution step. Users
 683 can check the fragment ion elution profiles from each inclusion list entry and how they are
 684 interconnected in the corresponding profile similarity network.



685

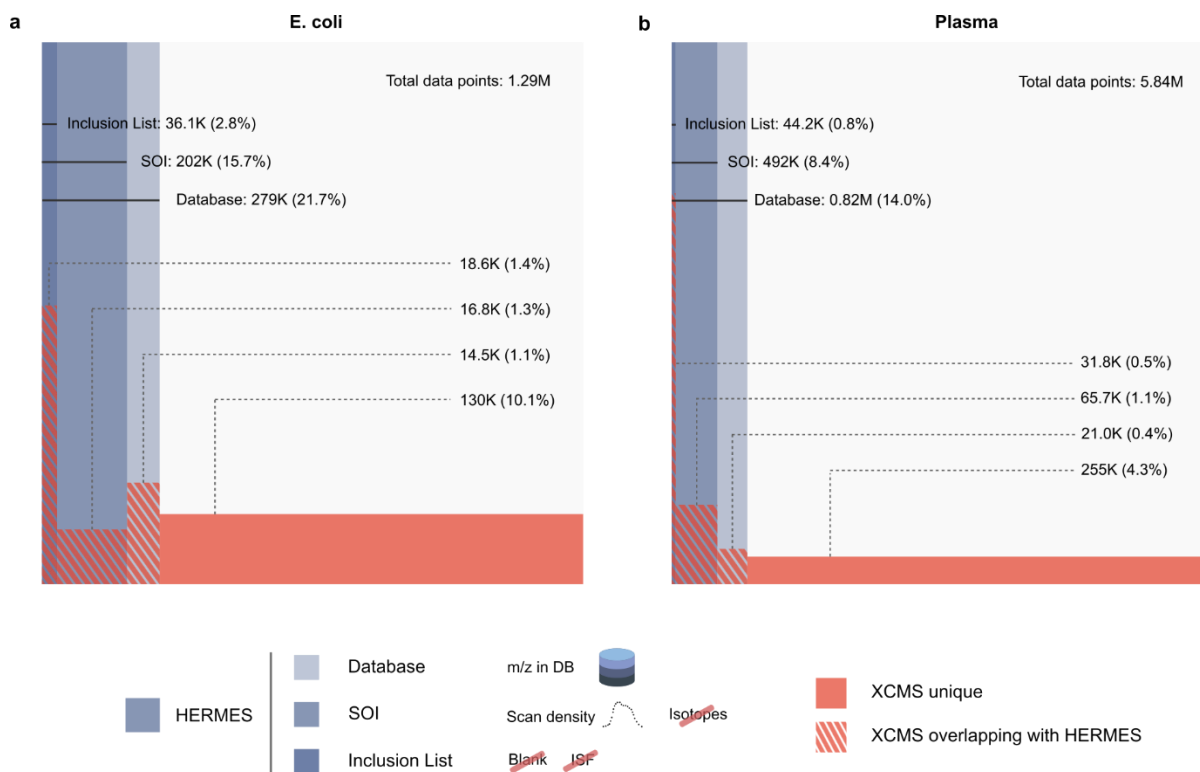
686

687 **Supplementary Figure S7. Discrimination of SOIs based on isotopic fidelity.** a) [M+H]⁺
 688 ion of chloridazon and b) [M+K]⁺ ion of 2-Amino-alpha-carboline overlapping at 0.27 ppm. The
 689 arrows indicate the characteristic [³⁷Cl] isotopologue present in chloridazon and the [⁴¹K]
 690 isotopologue absent in 2-Amino-alpha-carboline. The absence of characteristic isotopologue
 691 signals (Cl, Br, K, etc.) in an intense SOI results in a low isotopic fidelity score and its removal.



692

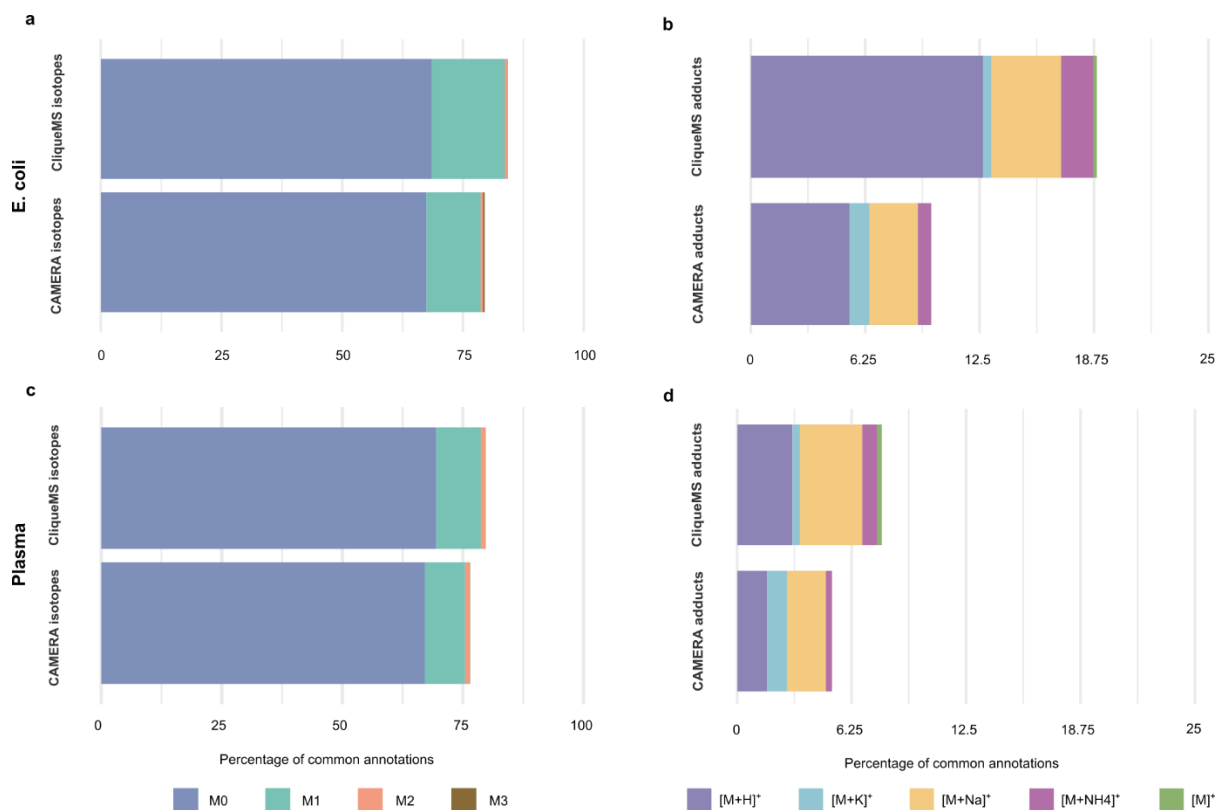
693 **Supplementary Figure S8.** Distribution of inclusion list entries of surface water in a) positive
694 and b) negative ionization mode after blank subtraction. The entries are coloured according to
695 their intensity values.



696

697

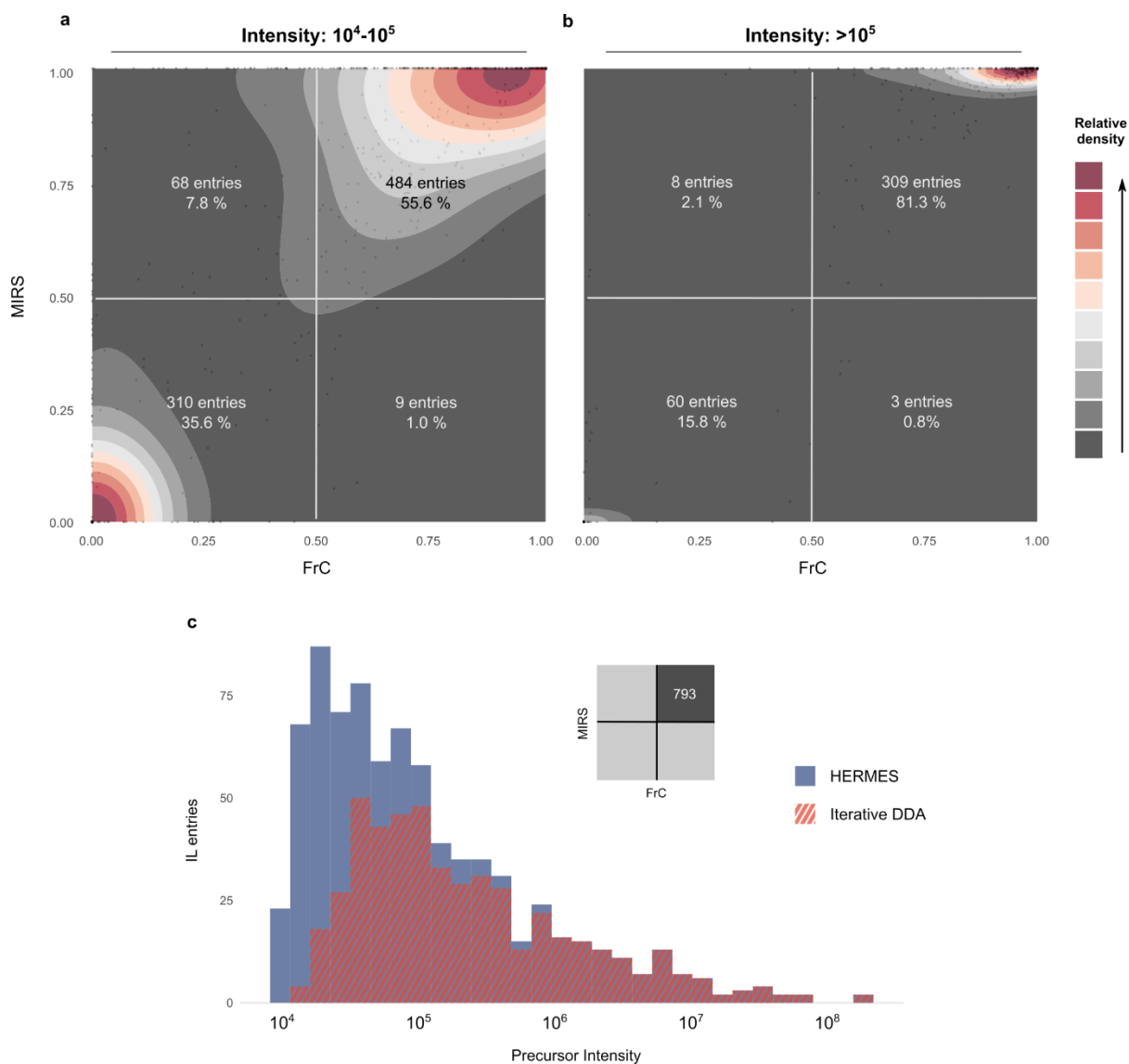
698 **Supplementary Figure S9. Venn-like diagram of the distribution of negative ionization**
 699 **LC/MS1 data points in different steps of the HERMES workflow and XCMS peak-**
 700 **associated points.** a) *E. coli* and b) human plasma extract. Database: data points that match
 701 any m/z from the ionic formula database (including isotopes). SOI: monoisotopic (M0)-
 702 annotated data points that are present in an unfiltered SOI list. Inclusion List: data points
 703 present in a filtered SOI list (including blank subtraction, isotopic filter and ISF removal steps).
 704 Percentages refer to the total number of LC/MS1 data points.



705

706

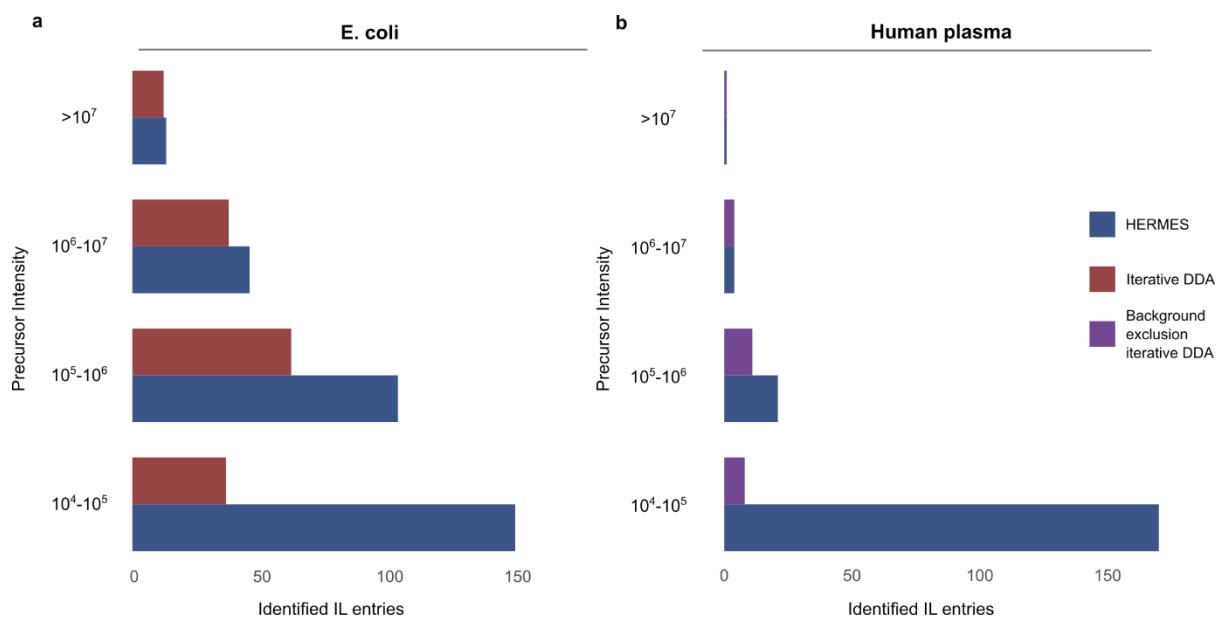
707 **Supplementary Figure S10: Comparing LC/MS1 annotation performance with CAMERA**
 708 **and CliqueMS.** Positive ionization data. a) and b) *E. coli*, c) and d) human plasma extract.
 709 Percentages refer to the set of datapoints annotated by HERMES and were detected as a
 710 peak by XCMS (see Online Methods for parameters used). The isotope annotation overlap (a
 711 and c) was high, due primarily to M0 annotations. On the other hand, adduct annotation
 712 overlap (b and d) was markedly low (<20% of datapoints matched the annotation).



713

714

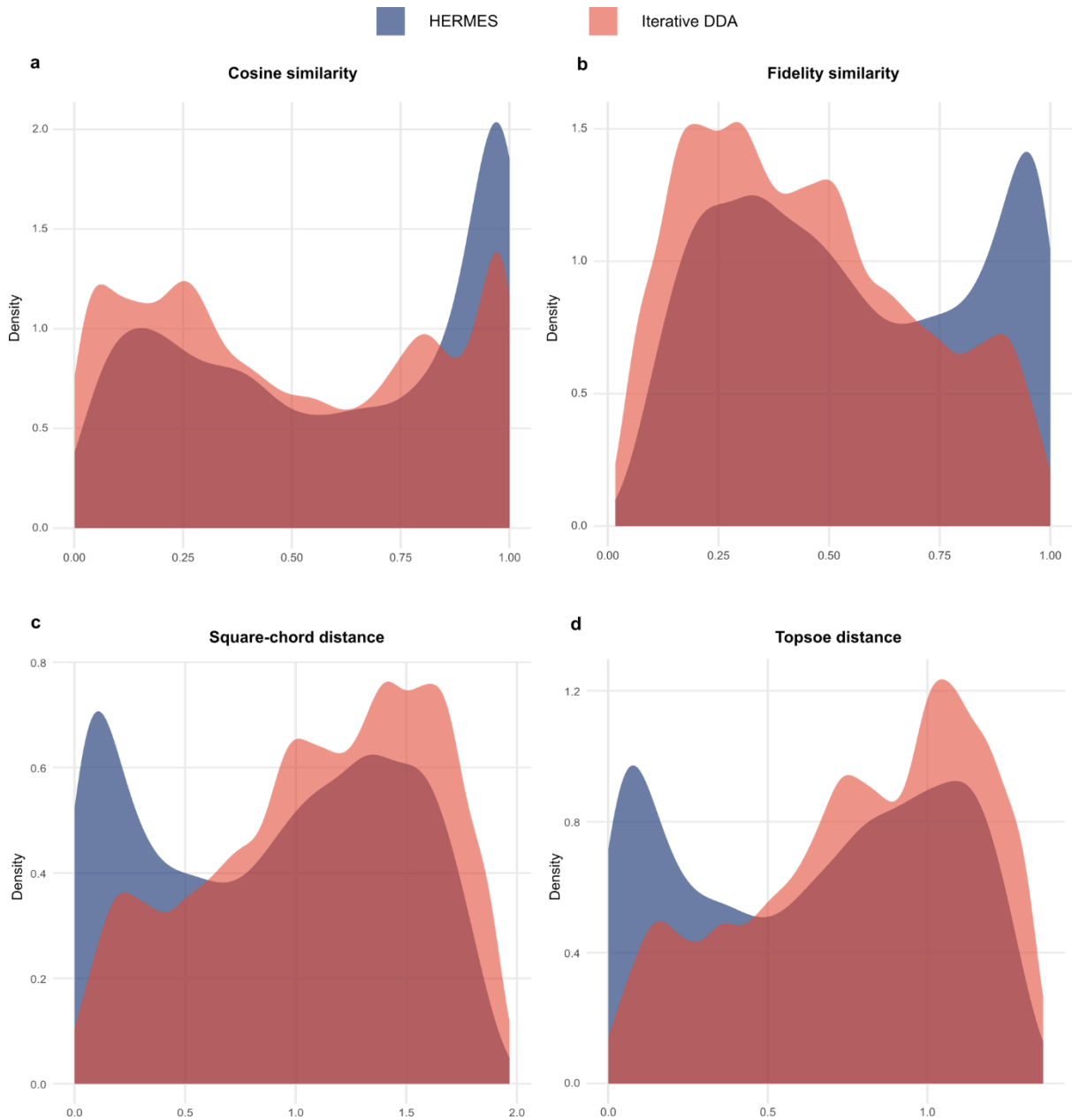
715 **Supplementary Figure S11. ^{13}C -enrichment distribution according to the precursor**
 716 **intensity.** a) and b) ^{13}C -enriched metabolites (FC and MIRS > 0.5) are mainly associated with
 717 abundant ions (intensity $>10^5$), while unlabeled precursors (FC and MIRS < 0.5) relate more
 718 frequently to low abundant ions (intensity between 10^4 - 10^5). c) ^{13}C -labeled precursors in
 719 iterative DDA corresponded to highly abundant ions that were also covered by Hermes.
 720 However, 56% of labelled low abundant ions were not covered by the iterative DDA.



721

722

723 **Supplementary Figure S12. Identified inclusion list entries according to the MS1**
 724 **precursor intensity in negative ionization data.** An inclusion list entry is considered
 725 identified if at least one MS2 scan associated with it has a compound hit in the reference MS2
 726 database with either cosine score > 0.8 (in-house database from MassBankEU, MoNA, Riken
 727 and NIST14 spectra), or Match > 90 and Confidence > 30 (mzCloud). a) *E. coli* extract. b)
 728 Human plasma extract.

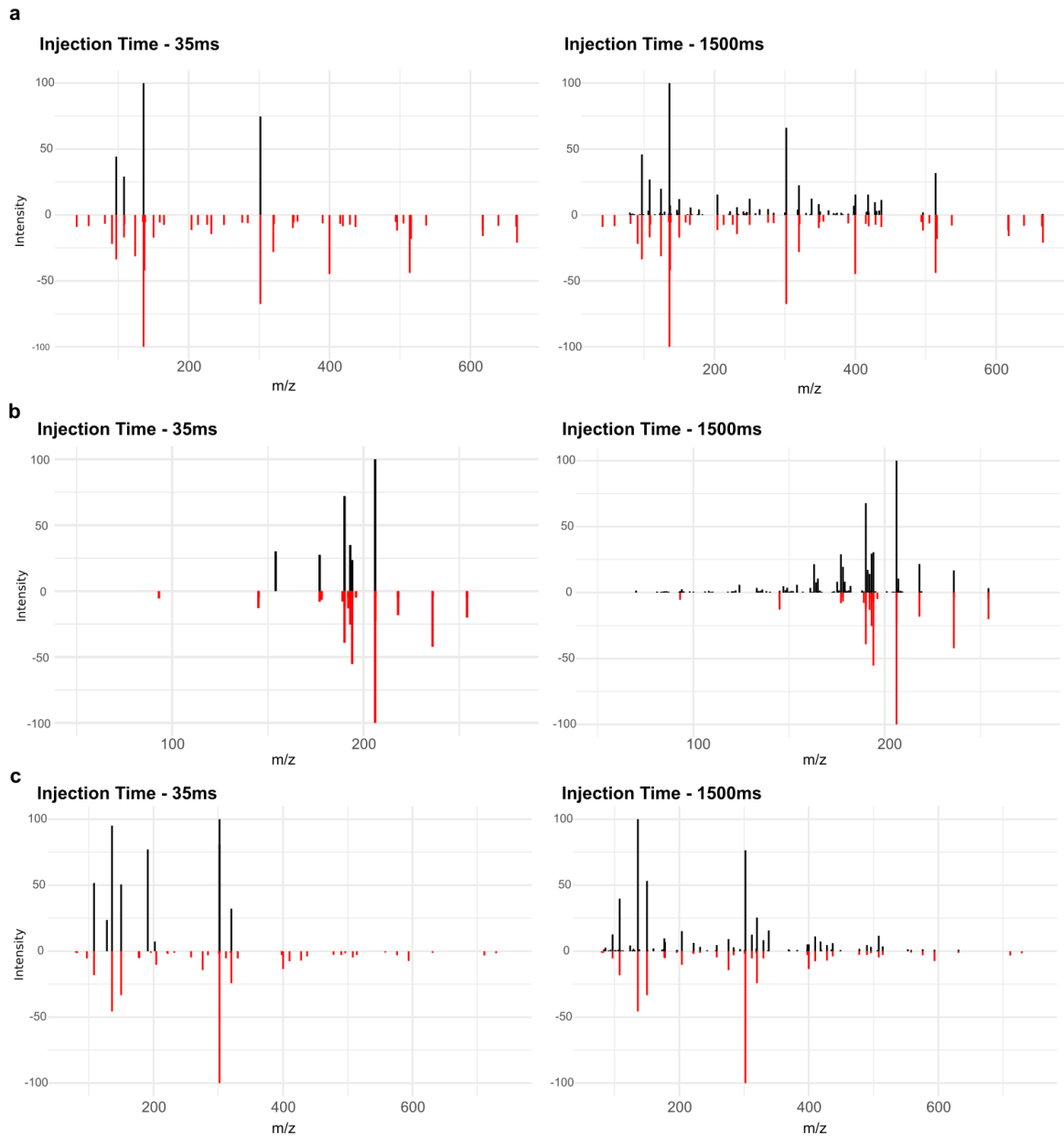


729

730

731 **Supplementary Figure S13. Alternative spectral similarity algorithms and spectrum-**
 732 **spectrum match scores.** a) Cosine similarity distribution b) Fidelity similarity distribution. c)
 733 Square-chord distance distribution. d) Topsoe distance distribution. A density estimation was
 734 calculated and normalized so that the integral of the curve equals 1. HERMES spectra showed
 735 higher similarity scores (a and b) and lower spectral distances (c and d) than DDA spectra.

736



737
738
739
740
741
742
743
744
745

Supplementary Figure S14. Injection time comparison (35 ms vs 1,500 ms). Experimental MS2 spectra (black) of a) NADH, b) Biopterin and c) NADPH against library spectra (red). All precursor ions had an intensity below 10^5 . A higher injection time resulted in richer spectra, with more matching fragments against the reference spectra and overall better matching scores.

746 **Supplementary Tables**747 **Supplementary Table 1.** Execution times and RAM usage up to the MS1 annotation step.
748

Number of files (Source)*	Total size of files (MB)	Number of unique ionic formulas	Used RAM (MB)*	Execution time#
2 (<i>E. coli</i>)	42.4	12,010	604	54s
4 (<i>E. coli</i>)	82.2	12,010	937	1 min 27s
8 (<i>E. coli</i>)	170	12,010	1,279	2 min 30s
2 (Plasma)	131	212,666	6,134	20 min 17s
4 (Plasma)	257	212,666	6,814	24 min 5s
8 (Plasma)	513	212,666	8,705	44 min 17s

749 *mzML files must be centroided. *Difference in RStudio RAM usage before and after
750 processing the files. #Until *Peaklist generation* (MS1 annotation) step.

751

752 **Supplementary Table 2.** List of spiked compounds.
753

Compound name	Formula	Adduct	m/z	RT (min)
1-(3,4-dichlorophenyl)-3-methylurea	C ₈ H ₈ Cl ₂ N ₂ O	[M+H] ⁺	219.0086	14.28
1-(3,4-Dichlorophenyl)-urea	C ₇ H ₆ Cl ₂ N ₂ O	[M+H] ⁺	204.9930	13.29
2,4-Dichloroaniline	C ₆ H ₅ Cl ₂ N	[M+H] ⁺	161.9872	16.77
2,6-dichlorobenzamide (BAM)	C ₇ H ₅ Cl ₂ NO	[M+H] ⁺	189.9821	8.18
2-aminoacetophenone	C ₈ H ₉ NO	[M+H] ⁺	136.0757	11.93
Atrazine	C ₈ H ₁₄ CIN ₅	[M+H] ⁺	216.1011	14.54
Azinphos-methyl	C ₁₀ H ₁₂ N ₃ O ₃ PS ₂	[M+H] ⁺	318.0131	17.17
Bezafibrate	C ₁₉ H ₂₀ CINO ₄	[M+H] ⁺	362.1154	15.95
Bromacil	C ₉ H ₁₃ BrN ₂ O ₂	[M+H] ⁺	261.0233	12.43
Caffeine	C ₈ H ₁₀ N ₄ O ₂	[M+H] ⁺	195.0877	6.83
Carbamazepine	C ₁₅ H ₁₂ N ₂ O	[M+H] ⁺	237.1022	13.27
Carbendazim	C ₉ H ₉ N ₃ O ₂	[M+H] ⁺	192.0768	6.38
Chlorpyrifos-ethyl	C ₉ H ₁₁ Cl ₃ NO ₃ PS	[M+H] ⁺	349.9336	23.34
Chlortoluron	C ₁₀ H ₁₃ CIN ₂ O	[M+H] ⁺	213.0789	14.31
Chloridazon	C ₁₀ H ₈ CIN ₃ O	[M+H] ⁺	222.0429	9.79
DEET	C ₁₂ H ₁₇ NO	[M+H] ⁺	192.1383	14.83
Desethylatrazine	C ₆ H ₁₀ CIN ₅	[M+H] ⁺	188.0698	9.78
Desisopropyl Atrazine	C ₅ H ₈ CIN ₅	[M+H] ⁺	174.0541	7.69

Diclofenac	C ₁₄ H ₁₁ Cl ₂ NO ₂	[M+H] ⁺	296.0240	18.37
Dimethenamid-p	C ₁₂ H ₁₈ CINO ₂ S	[M+H] ⁺	276.0820	17.37
Dimethoate	C ₅ H ₁₂ NO ₃ PS ₂	[M+H] ⁺	230.0069	10.29
Dimethomorph (isomer 1)	C ₂₁ H ₂₂ CINO ₄	[M+H] ⁺	388.1310	16.18
Dimethomorph (isomer 2)	C ₂₁ H ₂₂ CINO ₄	[M+H] ⁺	388.1310	16.59
Diuron	C ₉ H ₁₀ Cl ₂ N ₂ O	[M+H] ⁺	233.0243	15.07
Ethofumesate	C ₁₃ H ₁₈ O ₅ S	[M+H] ⁺	287.0948	18.48
Phenazone	C ₁₁ H ₁₂ N ₂ O	[M+H] ⁺	189.1022	8.66
Isoproturon	C ₁₂ H ₁₈ N ₂ O	[M+H] ⁺	207.1492	14.93
Linuron	C ₉ H ₁₀ Cl ₂ N ₂ O ₂	[M+H] ⁺	249.0192	17.24
Metazachlor	C ₁₄ H ₁₆ CIN ₃ O	[M+H] ⁺	278.1055	15.87
Metobromuron	C ₉ H ₁₁ BrN ₂ O ₂	[M+H] ⁺	259.0077	15.60
Metolachlor	C ₁₅ H ₂₂ CINO ₂	[M+H] ⁺	284.1412	18.92
Metoprolol	C ₁₅ H ₂₅ NO ₃	[M+H] ⁺	268.1907	9.46
Metoxuron	C ₁₀ H ₁₃ CIN ₂ O ₂	[M+H] ⁺	229.0738	11.94
Metribuzin	C ₈ H ₁₄ N ₄ OS	[M+H] ⁺	215.0961	13.19
Monuron	C ₉ H ₁₁ CIN ₂ O	[M+H] ⁺	199.0633	12.66
Nicosulfuron	C ₁₅ H ₁₈ N ₆ O ₆ S	[M+H] ⁺	411.1081	12.24
Pentoxifylline	C ₁₃ H ₁₈ N ₄ O ₃	[M+H] ⁺	279.1452	9.46
Pirimicarb	C ₁₁ H ₁₈ N ₄ O ₂	[M+H] ⁺	239.1503	9.11
Simazine	C ₇ H ₁₂ CIN ₅	[M+H] ⁺	202.0854	12.50
Sulfadimidine	C ₁₂ H ₁₄ N ₄ O ₂ S	[M+H] ⁺	279.0910	8.38
Sulfamethoxazole	C ₁₀ H ₁₁ N ₃ O ₃ S	[M+H] ⁺	254.0594	10.69
Terbutylazine	C ₉ H ₁₆ CIN ₅	[M+H] ⁺	230.1167	16.85
Tetraglyme	C ₁₀ H ₂₂ O ₅	[M+H] ⁺	223.1540	7.78
Triethyl Phosphate	C ₆ H ₁₅ O ₄ P	[M+H] ⁺	183.0781	10.94
Triphenylphosphine Oxide	C ₁₈ H ₁₅ OP	[M+H] ⁺	279.0933	15.34
Tri-n-butyl-phosphate	C ₁₂ H ₂₇ O ₄ P	[M+H] ⁺	267.1720	20.52
Tri-(2-chloroisopropyl)Phosphate	C ₉ H ₁₈ Cl ₃ O ₄ P	[M+H] ⁺	327.0081	17.24
Tris(2-chloroethyl)Phosphate (TCEP)	C ₆ H ₁₂ Cl ₃ O ₄ P	[M+H] ⁺	284.9612	14.26
2,4,6-trichlorophenol	C ₆ H ₃ Cl ₃ O	[M+H] ⁻	194.9177	17.77
2,4-dichlorophenol	C ₆ H ₄ Cl ₂ O	[M+H] ⁻	160.9566	16.52
2,4-dichlorophenoxyacetic Acid (2,4-D)	C ₈ H ₆ Cl ₂ O ₃	[M+H] ⁻	218.9621	15.25
2,4-dinitrophenol	C ₆ H ₄ N ₂ O ₅	[M+H] ⁻	183.0047	13.21

(4-chloro-2-methylphenoxy)-Acetic Acid (MCPA)	C ₉ H ₉ ClO ₃	[M+H] ⁻	199.0168	15.31
Bentazon	C ₁₀ H ₁₂ N ₂ O ₃ S	[M+H] ⁻	239.0496	14.44
Dichlorprop (2.4-DP)	C ₉ H ₈ Cl ₂ O ₃	[M+H] ⁻	232.9778	16.52
Mecoprop (MCP)	C ₁₀ H ₁₁ ClO ₃	[M+H] ⁻	213.0324	16.53
p,p-sulfonyldiphenol	C ₁₂ H ₁₀ O ₄ S	[M+H] ⁻	249.0227	11.21
N-acetyl sulfamethoxazole	C ₁₂ H ₁₃ N ₃ O ₄ S	[M+H] ⁺	296.0700	11.06
Metolachlor ESA	C ₁₅ H ₂₃ NO ₅ S	[M+H] ⁺	330.1370	11.24
10,11-dihydro-10,11-dihydroxy Carbamazepine	C ₁₅ H ₁₄ N ₂ O ₃	[M+H] ⁺	271.1077	7.55
Gabapentin	C ₉ H ₁₇ NO ₂	[M+H] ⁺	172.1332	6.45
Hydrochlorothiazide	C ₇ H ₈ ClN ₃ O ₄ S ₂	[M+H] ⁻	295.9572	7.20
Desfenylchloridazon	C ₄ H ₄ ClN ₃ O	[M+H] ⁺	146.0116	2.25
Lamotrigine	C ₉ H ₇ Cl ₂ N ₅	[M+H] ⁺	256.0151	9.36
Metazachlor ESA	C ₁₄ H ₁₇ N ₃ O ₄ S	[M+H] ⁺	324.1013	9.19
N-formyl-4-aminoantipyrine	C ₁₂ H ₁₃ N ₃ O ₂	[M+H] ⁺	232.1081	7.12
N-acetyl-4-aminoantipyrine	C ₁₃ H ₁₅ N ₃ O ₂	[M+H] ⁺	246.1237	7.08
Metazachlor OA	C ₁₄ H ₁₅ N ₃ O ₃	[M+H] ⁺	274.1186	9.32
Sitagliptin	C ₁₆ H ₁₅ F ₆ N ₅ O	[M+H] ⁺	408.1254	10.27
Valsartan Acid	C ₁₄ H ₁₀ N ₄ O ₂	[M+H] ⁺	267.0877	11.78
Gabapentin-lactam	C ₉ H ₁₅ NO	[M+H] ⁺	154.1226	11.22
HMMM	C ₁₅ H ₃₀ N ₆ O ₆	[M+H] ⁺	391.2300	13.24
Candesartan	C ₂₄ H ₂₀ N ₆ O ₃	[M+H] ⁺	441.1670	14.37
Irbesartan	C ₂₅ H ₂₈ N ₆ O	[M+H] ⁺	429.2397	14.13
Valsartan	C ₂₄ H ₂₉ N ₅ O ₃	[M+H] ⁺	436.2343	16.51
Sebutylazine	C ₉ H ₁₆ ClN ₅	[M+H] ⁺	230.1167	16.20
Telmisartan	C ₃₃ H ₃₀ N ₄ O ₂	[M+H] ⁺	515.2442	14.07
Cetirizine	C ₂₁ H ₂₅ ClN ₂ O ₃	[M+H] ⁺	389.1627	23.33
1-H-benzotriazole	C ₆ H ₅ N ₃	[M+H] ⁺	120.0556	7.92
4-methyl-1H-benzotriazole	C ₇ H ₇ N ₃	[M+H] ⁺	134.0713	9.94
5-methyl-1H-benzotriazole	C ₇ H ₇ N ₃	[M+H] ⁺	134.0713	10.07
5,6-dimethyl-1H-benzotriazole	C ₈ H ₉ N ₃	[M+H] ⁺	148.0869	11.52
5-chloro-1H-benzotriazole	C ₆ H ₄ ClN ₃	[M+H] ⁺	154.0167	11.30
2-aminobenzothiazole	C ₇ H ₆ N ₂ S	[M+H] ⁺	151.0324	6.43
2-hydroxybenzothiazole	C ₇ H ₅ NOS	[M+H] ⁺	152.0165	11.59

2-(methylthio)benzothiazole	C ₈ H ₇ NS ₂	[M+H] ⁺	182.0093	17.38
-----------------------------	---	--------------------	----------	-------

754 **Supplementary Material: Algorithms**

755 **Algorithm 1:** Resolution-adapted isotopic envelope calculation

756 **Input:** Ionic formula f , corresponding m/z , instrument resolution R at $m/z = 200$, separation
757 threshold k

758 **Output:** List of resolvable isotopes I

759 $R' = R(m/z)$ // Estimate instrument resolution for the ionic formula m/z

760 $d = 2 \cdot k \cdot R'$ // Calculate minimum resolvable m/z distance d

761 $I_{th} = iso_pattern(f)$ // Calculate theoretical isotopic pattern

762 $N_{cluster} =$ Number of isotope clusters in I_{th} such that $dist(i, j) < 0.1 Da, \forall i, j \in Cluster$

763 For $cluster$ from 1 to $N_{cluster}$:

764 $iso_{mz} = \text{sort}(mz_i) \forall i \in cluster$

765 if any ($iso_{mz} > d$): // If there's more than one distinguishable isotope

766 Split cluster into groups G such that: $\forall i, j \in G \ dist(i, j) < d$

767 For each $group$ in G :

768 Add the isotope with highest intensity in $group$ to I

769 Else:

770 Add the isotope with highest intensity in $cluster$ to I

771

772 Return I

773

774 **Algorithm 1:** By considering the instrumental resolution, this algorithm can adaptively
775 annotate isotopes according to the MS1 acquisition settings. The idea is to extend the
776 definition of chromatographic peak resolution (which compares peaks RTs and widths) to the
777 m/z profiles of isotope peaks and infer whether the centroidization algorithm (frequently used
778 in mzML file pre-processing) can split them into separate data points (see Supplementary
779 Figure 1). To do this, the algorithm (a) estimates the instrument *local* resolution for a particular
780 ionic formula based on a reference resolution, (b) calculates the minimum distinguishable
781 distance d for the centroidization algorithm, (c) calculates a theoretical isotopic envelope and
782 (d) finds all groups of isotopes with distances $< d$.

783

784

785

786 **Algorithm 2:** Scans of Interest (SOI) detection

787 **Input:** Annotated data points matrix $D = \{RT, intensity\}$, retention time bin width bw ,
788 minimum intensity I_{min} , acquired scans header h , minimum scan density (%) ρ , minimum
789 chaos score C_{min}

790 **Output:** List of SOIs SOI_{list}

791 $D' = filter(D, intensity > I_{min})$

792 // Set RT bins

793 $B_{start} = [0, bw, 2 \cdot bw, \dots]$; $B_{end} = [bw, 2 \cdot bw, 3 \cdot bw, \dots]$

794 // Count number of acquired scans in each RT bin

795 $Thr = [B_{start} \leq h_{RT} \leq B_{end}] \cdot \rho$

796 For each bin :

797 $N_{bin} = B_{start}[bin] \leq D'_{RT} \leq B_{end}[bin]$ // Count annotated data points in the bin

798 $\phi_{bin} = N_{bin} > Thr_{bin}$ // Determine whether there are enough data points

799 $\sigma = S(\phi_{bin})$ // Find all 1D-connected subsets in ϕ_{bin}

800 For each $element$ in σ :

801 Calculate ρ_{chaos} with the $element$ intensity vector (see Algorithm 4)

802 If $\rho_{chaos} > C_{min}$:

803 Add $element$ information to the SOI_{list} (start, end, intensity, annotation ...)

804 Return SOI_{list}

805

806 **Algorithm 2:** The SOI detection algorithm establishes a new way of detecting peaks without
807 imposing a gaussian peak-shape on the data. By using a moving window to calculate an
808 average scan density, the algorithm is robust against missing data points (which are frequent
809 in low-intensity signals). SOIs are detected from an already annotated set of datapoints –
810 defined by the considered formulas and adducts– so the resulting SOIs bypass the limitations
811 of current annotation methods (mainly, the requirement of adduct and/or isotope peaks to
812 annotate).

813

814 **Algorithm 3:** Continuous MS2 scan deconvolution

815 **Input:** Continuously acquired MS2 scans S from a common precursor m/z , instrument ppm
816 error ppm

817 **Output:** Deconvoluted MS2 spectra

818 **// Mass trace detection**

819 For each $scan$ in S :

820 Match the scan m/z to those of a trace list T

821 If there's a match:

822 Extend the trace

823 If there are no matches:

824 Start a new trace and add the scan

825 Split traces with RT gaps $> 3s$

826 Run *Centwave* peak detection algorithm on each individual trace

827 **// Information propagation based on elution profile similarity**

828 For each trace t detected as a peak:

829 For each trace t' **not** detected as a peak:

830 Calculate the Pearson correlation r between t and t'

831 If $r > 0.8$:

832 Trim t' RT interval to match t

833 Consider t' as a peak from now on

834 Repeat previous loop until no more peaks are detected

835 **// Network analysis**

836 Calculate cosine similarities between each trace, generating a similarity matrix M

837 Generate a network N from M

838 Find all connected subgraphs in N

839 For each subgraph:

840 Apply a greedy partitioning algorithm.

841 If the modularity resulting from the partition is $> x$:

842 Split the subgraph according to the partition

843 For each resulting subgraph:

844 Retrieve the traces that form the subgraph

845 Assemble a MS2 spectrum from the maximum intensities of each trace

846 Return a list of the assembled MS2 spectra.

847

848 **Algorithm 3:** The MS2 deconvolution algorithm separates fragments from coeluting
849 compounds with a common precursor m/z . A key feature of the algorithm is the use of
850 *Centwave* peak-picking and the *information propagation* from one fragment mass trace to
851 another applying a Pearson correlation: if a peak is detected with *Centwave* in one trace
852 (usually the most intense), similarly-shaped, coeluting traces are treated as peaks
853 independently of the peak-picking. This propagation step is particularly useful when most
854 traces are weak and noisy, as they are frequently missed by peak-picking algorithms.

855 **Algorithm 4:** 1-dimensional ρ_{chaos} calculation

856 **Input:** Vector I of intensities of the data points corresponding to a SOI, number n of intensity
857 thresholds

858 **Output:** ρ_{chaos}

859 $I = \frac{I - \min(I)}{\max(I) - \min(I)}$ // Scale intensity to [0, 1]

860 For $level$ in 1 to n :

861 $\tau = level / n$

862 $\phi = I > \tau$ // Find all intensities above threshold

863 $\phi' = \psi(\phi)$ // Fill gaps of length 1

864 $N_{level_i} = S(\phi')$ // Find all 1D-connected subsets in ϕ'

865 $\rho_{chaos} = 1 - \frac{\sum_{i=1}^n N_{level_i}}{length(I) \cdot n}$

866 Return ρ_{chaos}

867

868 **Algorithm 4:** The ρ_{chaos} calculated by this algorithm is the 1D equivalent of the ρ_{chaos} applied
869 in MS-imaging software like METASPACE. The objective is to quantify the amount of structure
870 contained in a datapoint intensity vector: when increasing a threshold $level$, a vector of noisy
871 datapoints generates many unconnected sets above the $level$ (low ρ_{chaos}), while well-
872 structured datapoints stay in one/few connected sets (high ρ_{chaos}).

873