

1 **Metabolome Identification by Systematic Stable Isotope Labeling Experiments and False**
2 **Discovery Analysis with a Target-Decoy Strategy**

3
4 *Drew R. Jones^{1#}, Xusheng Wang^{2#}, Tim Shaw^{2,3}, Ji-Hoon Cho², Ping-Chung Chen^{1,4}, Suiping*
5 *Zhou², Yuxin Li, Nam Chul Kim⁵, J. Paul Taylor^{5,6}, Udhghatri Kolli⁷, Jiaxu Li⁷, and Junmin*
6 *Peng^{1,2,4*}*

7 ¹Department of Structural Biology, ²St. Jude Proteomics Facility, ³Department of Computational
8 Biology, ⁴Department of Developmental Neurobiology, ⁵Department of Cell and Molecular
9 Biology, ⁶Howard Hughes Medical Institute. ^{1,2,3,4,5}St. Jude Children's Research Hospital, 262
10 Danny Thomas Place, Memphis, TN 38105, USA

11 ⁷Department of Biochemistry, Molecular Biology, Entomology, and Plant Pathology, Mississippi
12 State University, 32 Creelman Street Mississippi State, MS 39762, USA

13 #These authors contributed equally to the work

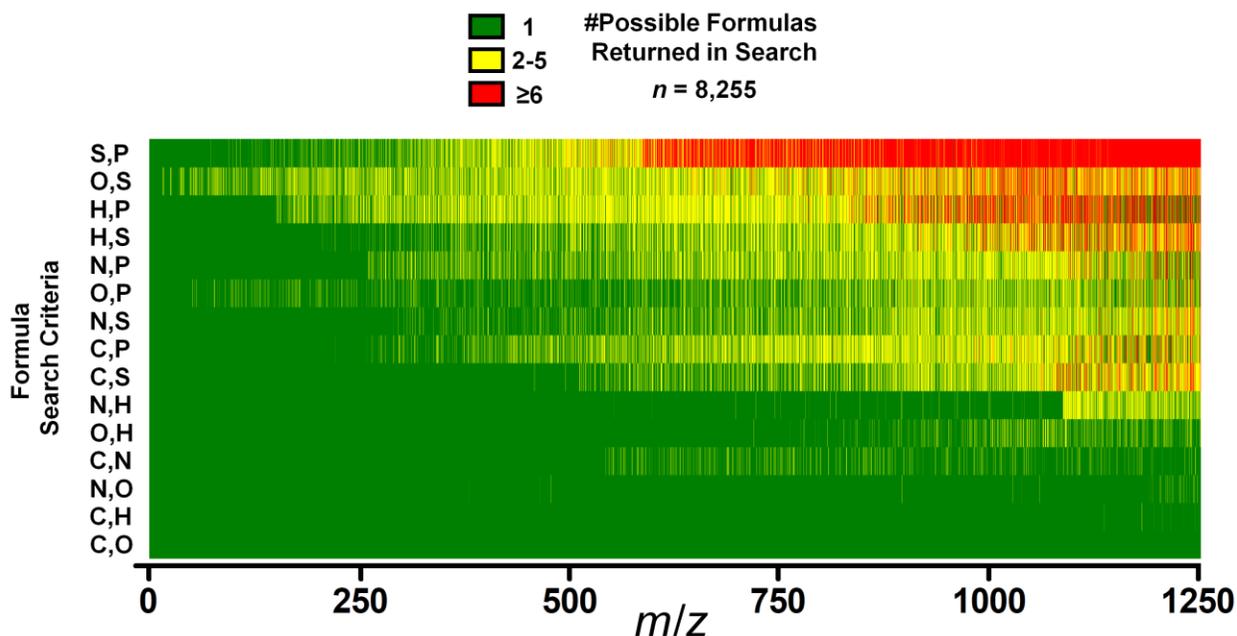
14 *Correspondence: Junmin Peng (junmin.peng@stjude.org)

15
16
17
18
19
20

Table of Contents

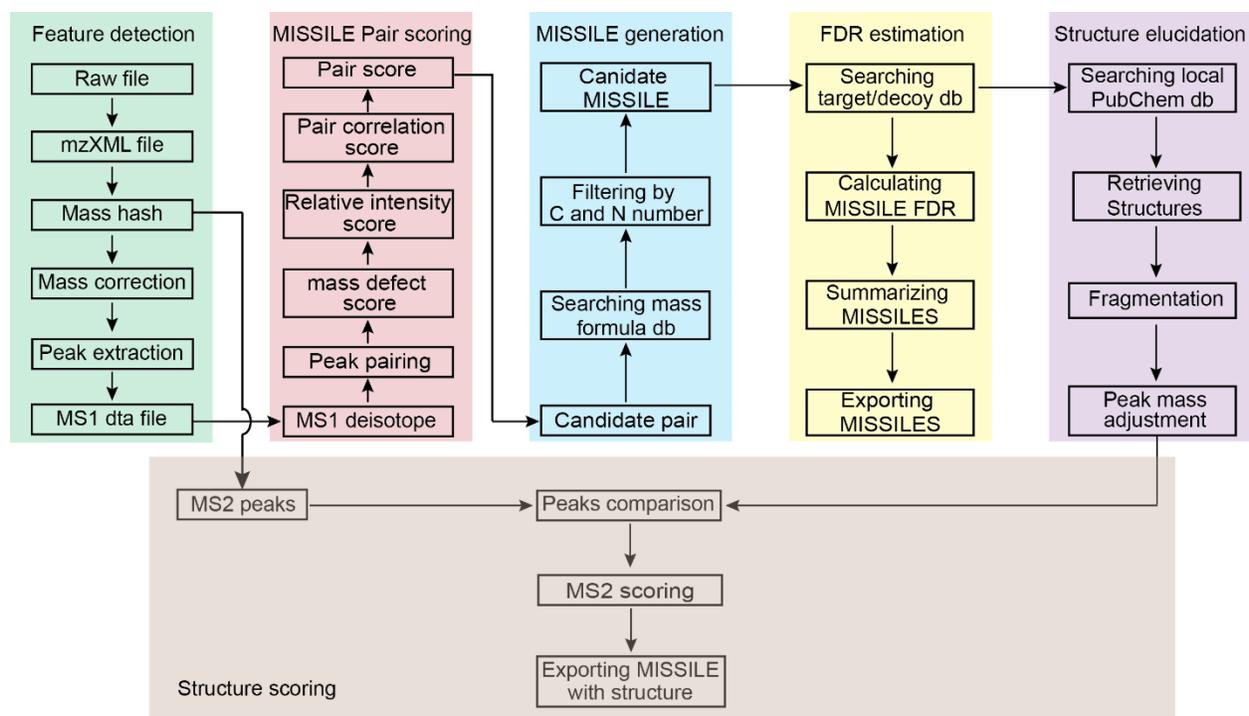
21 Supplementary Fig. 1. Heat map visualization of chemical formula searches.....3
22 Supplementary Fig. 2. Detailed workflow for JUMPm.....4
23 Supplementary Fig. 3. Distribution of signal and noise levels in the testing dataset.....5
24 Supplementary Fig. 4. Mass correction by the moving average method.....6
25 Supplementary Fig. 5. Other types of sample labeling compatible with JUMPm.....7
26 Supplementary Fig. 6. Conceptual workflow for JUMPm search of unlabeled and partially
27 labeled metabolites.....8
28 Supplementary Fig. 7. Determination of Pscore parameter tolerance.....9
29 Supplementary Fig. 8. Systematic parameter optimization.....10
30 Supplementary Fig. 9. Number of molecular formulas containing 16 selected elements in the
31 HMDB and YMDB.....11
32 Supplementary Fig. 10. Comparison of the number of formulas (a) and structures (b) between
33 the PubChem, HMDB, and YMDB databases.....12
34 Supplementary Fig. 11. Scoring of metabolite-spectrum matches.....13
35 Supplementary Fig. 12. FDR characterization with labeled yeast.....14

36 Supplementary Fig. 13. Stable isotope labeling in yeast.....15
 37 Supplementary Fig. 14. Purity of isotope incorporation in yeast metabolites.....16
 38 Supplementary Fig. 15. MS/MS annotation using stable isotope information.....17
 39 Supplementary Fig. 16. Analysis of a formic acid adduct from yeast.....18
 40 Supplementary Fig. 17. Spectral library search of standards and yeast MS2 scans.....19
 41

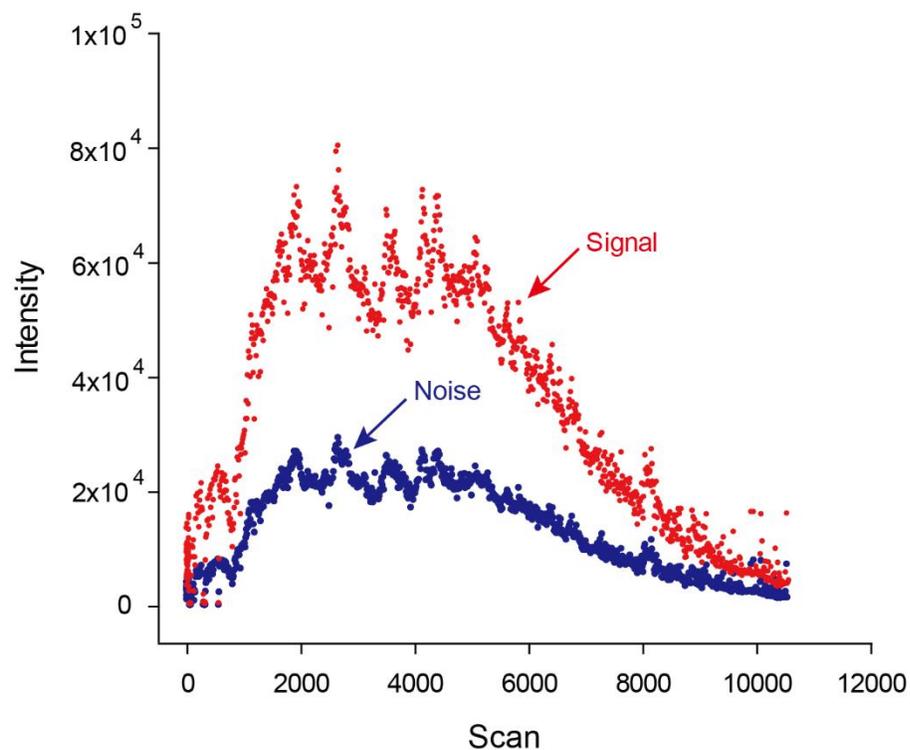


42
 43 **Supplementary Fig. 1. Heat map visualization of chemical formula searches restricted by**
 44 **known formula information (two elements) as indicated.** Each search was performed with a
 45 mass tolerance of 4 ppm. E.g. for row “C,N” the available information was the accurate
 46 theoretical mass (+/- 4ppm) of the HMDB entry and the number of carbon and nitrogen atoms in
 47 that entry’s chemical formula. Total number of searches = 8,255 for each row.

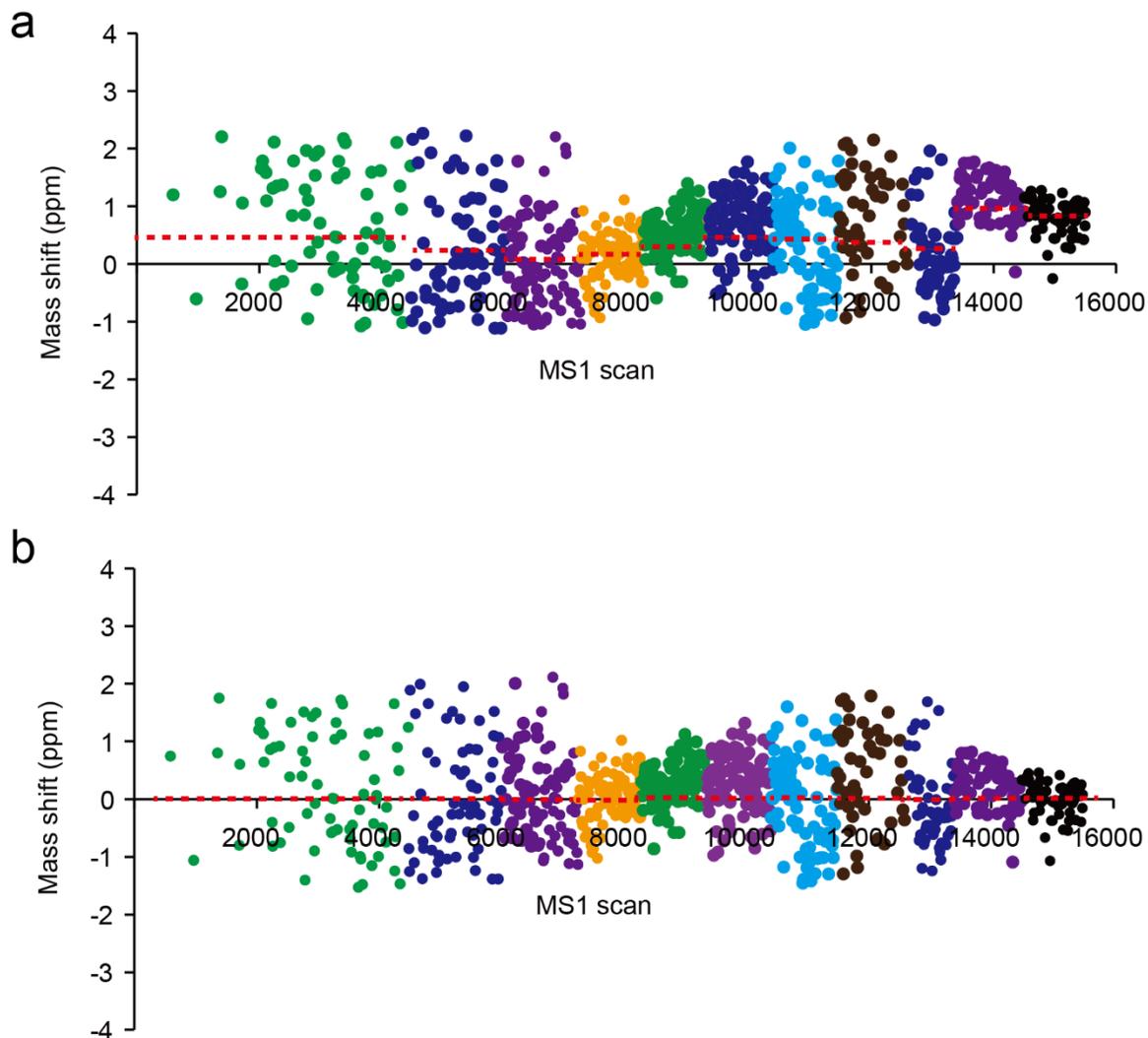
48
 49
 50
 51
 52
 53
 54
 55
 56
 57
 58
 59
 60
 61
 62



63
 64 **Supplementary Fig. 2. Detailed workflow for JUMPm.** JUMPm has six major components,
 65 including feature detection, isotope peak pairing, formula identification, FDR estimation,
 66 structural elucidation, and spectrum match scoring. JUMPm accepts .raw or .mzXML files
 67 (including MS1 and MS2) and outputs tables of identified metabolites and formulas. db =
 68 database, FDR = false discovery rate, MISSILE = chemical formula determined from accurate
 69 mass and labeled mass shift information.

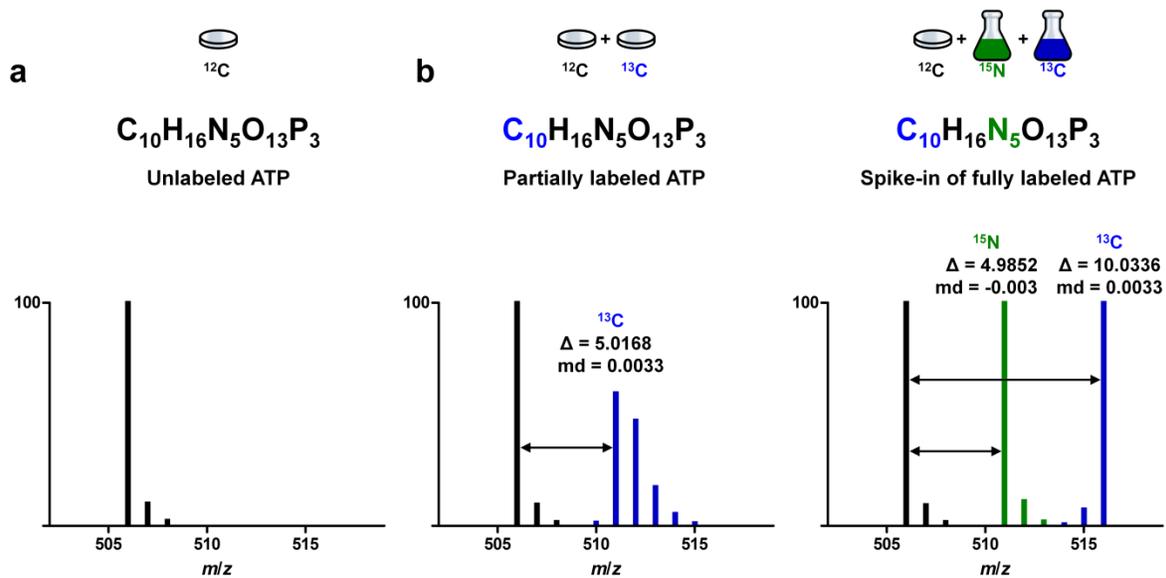


70
71 **Supplementary Fig. 3. Distribution of signal and noise levels in the testing dataset.**
72 JUMPm first defines noise peaks as those which cannot be repeatedly detected in adjacent
73 MS1 scans (i.e. peaks detected in only a single scan, **Online Methods**). For each scan, the
74 program collects these noise peaks, removes outliers and obtains the average intensity of the
75 remaining peaks. The average is used as scan-specific noise level to enable the calculation of
76 the signal-to-noise (S/N) ratio of all peaks. An S/N ratio of 3 is set as the default cutoff for the
77 peaks but is user-adjustable.
78
79



80
 81 **Supplementary Fig. 4. Mass correction by the moving average method.** (a) Mass accuracy
 82 can drift during the course of a single run, and therefore the use of moving windows enables
 83 accurate mass correction. The mass error is determined by the mass difference between the
 84 detected and theoretical mass of the polysiloxane ion (445.120025, **Online Methods**). The
 85 entire chromatogram is typically divided into 10 windows, each containing the same number of
 86 scans. The mean of the mass errors in each window (dashed red line) is used to perform mass
 87 correction (b) on the ions within that window.

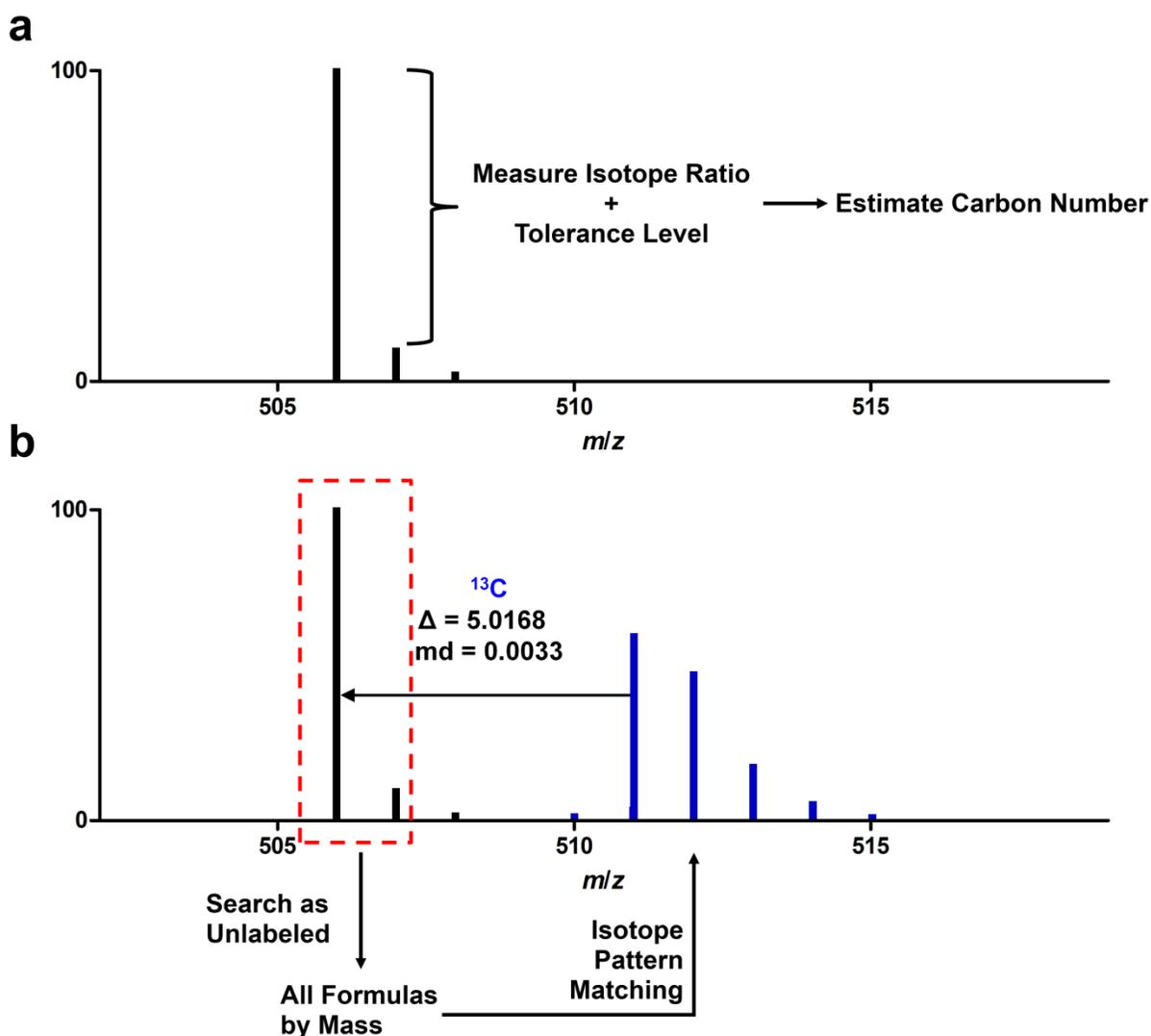
88
 89



90
91

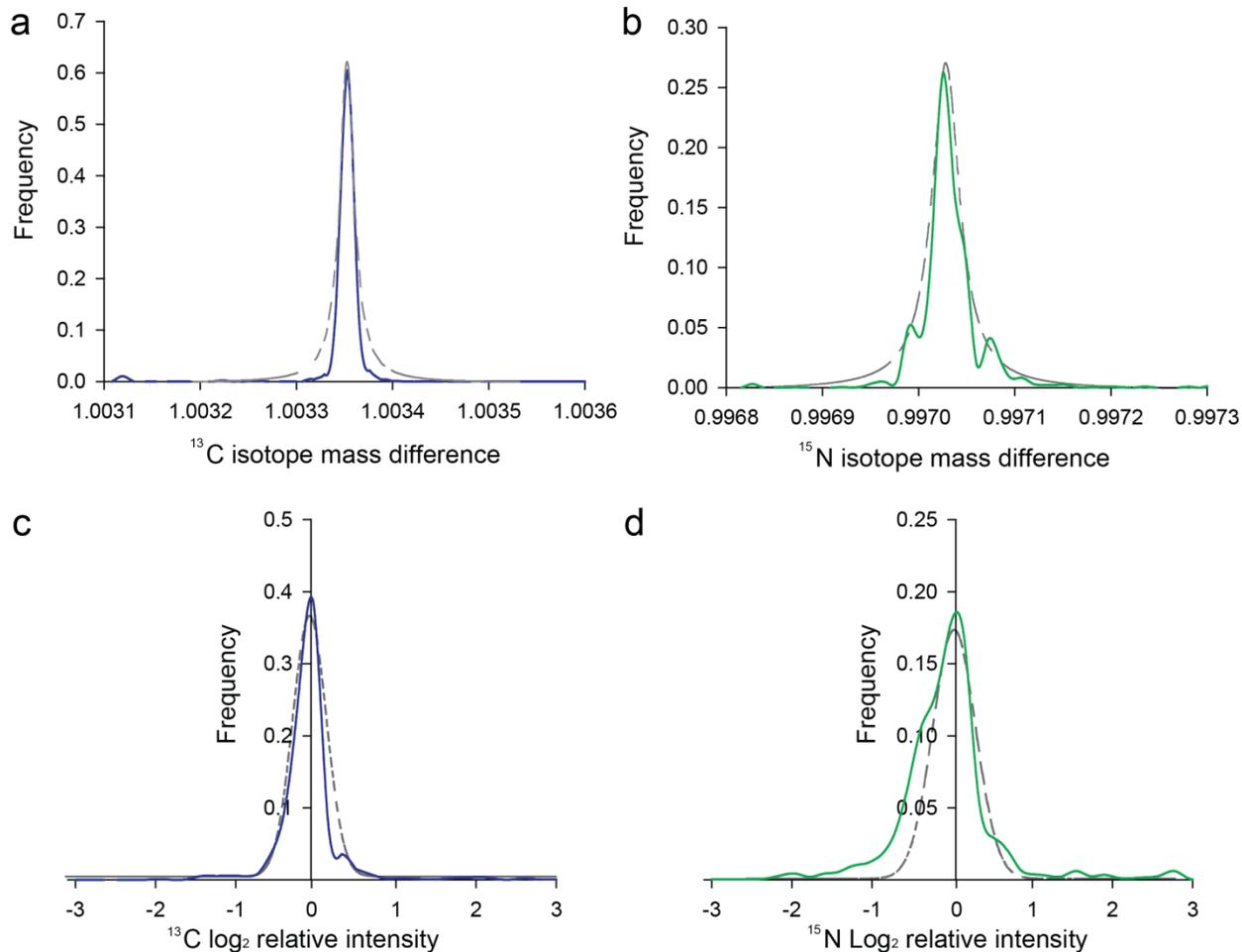
92 **Supplementary Fig. 5. Other types of sample labeling compatible with JUMPm.** (a)
 93 Example MS1 spectrum of unlabeled ($^{12}\text{C}^{14}\text{N}$) Adenosine triphosphate (ATP) from HEK293 cells
 94 (dish). Natural carbon atoms are ~98.9% pure, with ~1.1% of the atoms being the stable
 95 isotope carbon-13 (^{13}C). The tallest peak at 506 m/z is the monoisotopic peak, corresponding to
 96 ATP molecules composed of only ^{12}C carbon atoms. With 10 carbon atoms in the formula,
 97 there is a ~9% chance of any given ATP molecule containing 1 natural ^{13}C atom in its structure.
 98 Therefore, the peak at 507 m/z is ~9% as abundant as the monoisotopic peak. (b) an MS1
 99 spectrum from a mixture of unlabeled ATP and partially labeled ATP (blue) from HEK293 cells.
 100 Unlabeled and partially labeled cells were grown separately then mixed (dishes). Although ATP
 101 has 10 carbons in its formula, we observed a mass shift of 5 m/z to the tallest labeled peak.
 102 The purity of the labeling determines the relative intensity and pattern of the labeled peaks. Δ =
 103 mass shift, due to isotope label, md = mass defect, unique to each stable isotope. (c) MS1
 104 spectrum of ATP from a spike-in sample containing unlabeled ($^{12}\text{C}^{14}\text{N}$) HEK293 cells, $^{12}\text{C}^{15}\text{N}$
 105 yeast, and $^{13}\text{C}^{14}\text{N}$ yeast (beakers). The $^{12}\text{C}^{15}\text{N}$ and $^{13}\text{C}^{14}\text{N}$ yeast samples are fully labeled,
 106 showing a shift of 10 carbons and 5 nitrogens, in agreement with the chemical formula of ATP.
 107 Compound identity was confirmed by authentic standards.

108



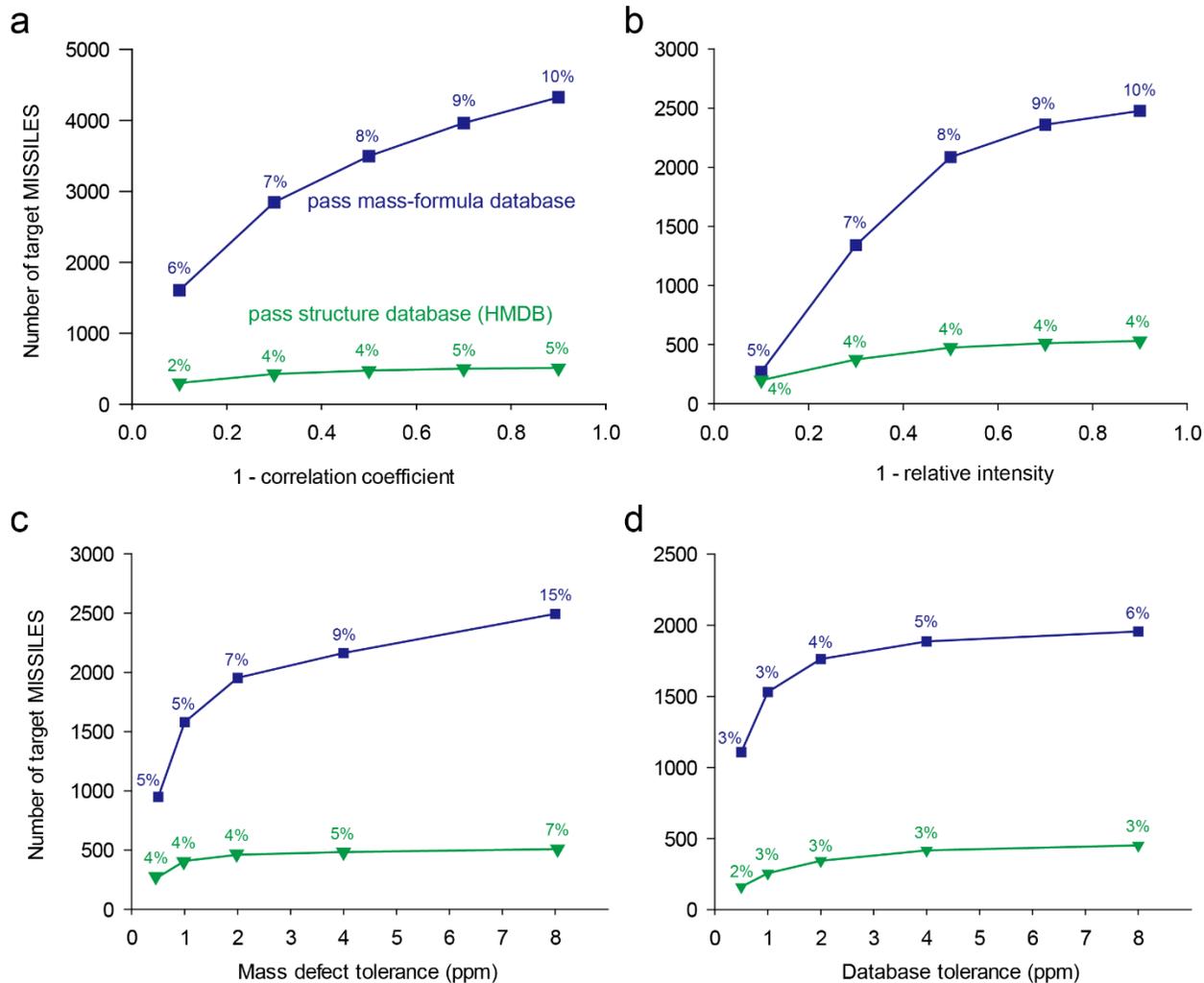
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123

Supplementary Fig. 6. Conceptual workflow for JUMPm search of unlabeled and partially labeled metabolites. (a) MS1 spectrum of ATP from unlabeled ($^{12}\text{C}^{14}\text{N}$) HEK293 cells alone. Compound identity was confirmed by authentic standards. The relative intensity of the A+1 (nature abundance isotope peak) is used to estimate the number of carbon atoms in the formula. (b) MS1 spectrum of unlabeled and partially labeled ATP from a mixture of unlabeled and $^{13}\text{C}^{14}\text{N}$ labeled HEK293 cells. First, JUMPm detects ion pairs with a mass defect of 0.0033, specific to the ^{13}C label. Then, JUMPm searches the unlabeled peak, using the A+1 intensity to estimate the carbon number and identify potential formulas. From this list of formulas, JUMPm uses isotope pattern simulation on the partially labeled $^{13}\text{C}^{14}\text{N}$ (blue) cluster to determine the best matching formula.



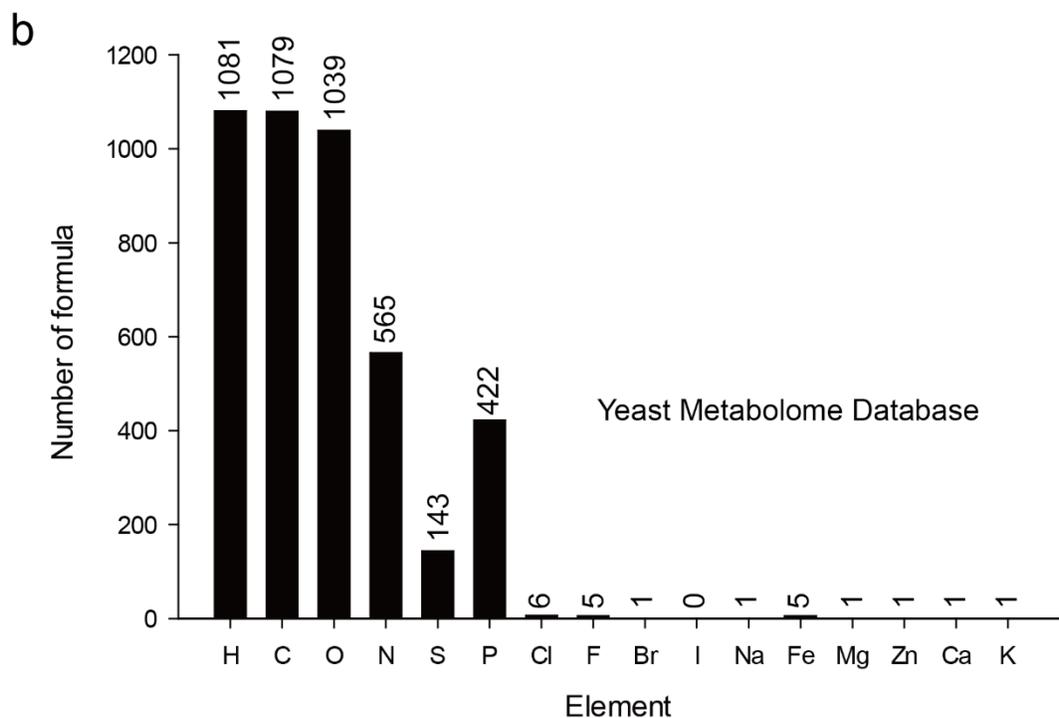
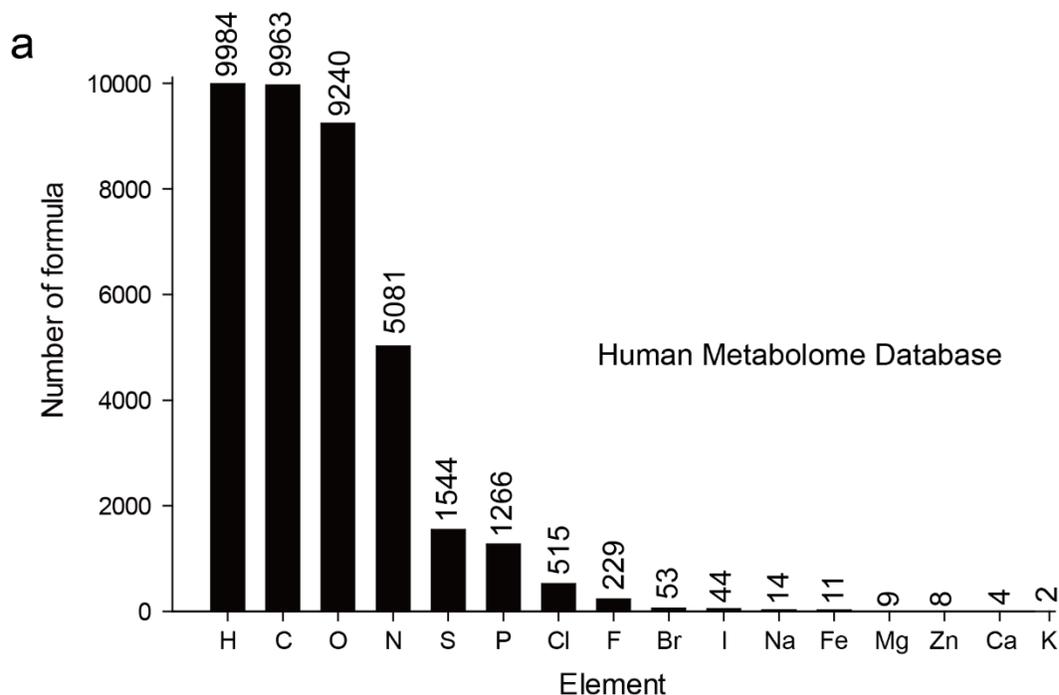
124
 125
 126
 127
 128
 129
 130
 131
 132
 133
 134
 135
 136
 137
 138
 139
 140
 141
 142

Supplementary Fig. 7. Determination of Pscore parameter tolerance. (a-b) Global frequency histograms for the isotope mass difference and relative intensity of the labeled peaks. The global tolerances for the ^{13}C and ^{15}N mass defects relative to the unlabeled monoisotopic peak were calculated by examining a population of the top 10% high confidence MISSILE formulas and their pairings. The mass defect is normalized about the monoisotopic mass, and is therefore primarily affected by mass resolution and mass accuracy, not the absolute mass value (small vs. large ions). The mass defect for both isotope labels was tightly distributed around the theoretical value. The nitrogen histogram is typically noisier because there are fewer nitrogen atoms in metabolites compared to carbon. This effectively reduces the sample size. (c-d) Relative ion intensity between labels also showed a normal distribution around $\pm 20\%$. Systematic shifts reflect unequal mixing of the samples while biological variation in metabolite levels contributes to the variation. For quantitative studies, the global ratio of labels can be used to fine-tune the mixture of samples for better quantitative sensitivity.

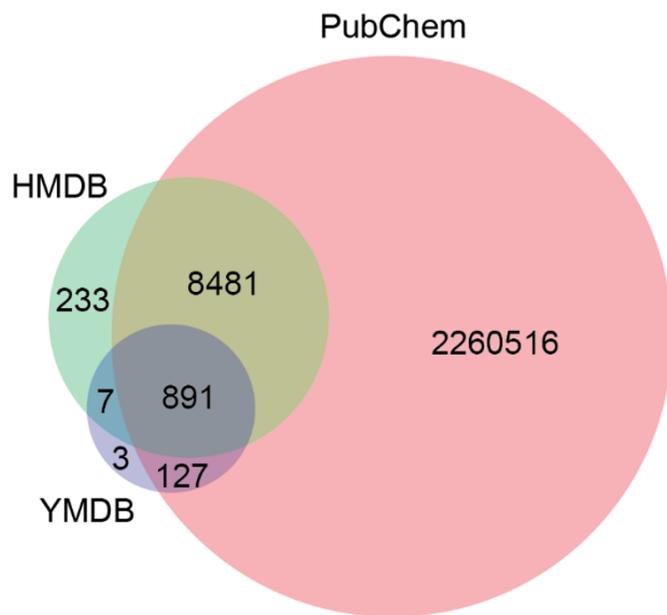


143
 144
 145
 146
 147
 148
 149
 150
 151
 152
 153
 154
 155
 156
 157
 158
 159
 160
 161

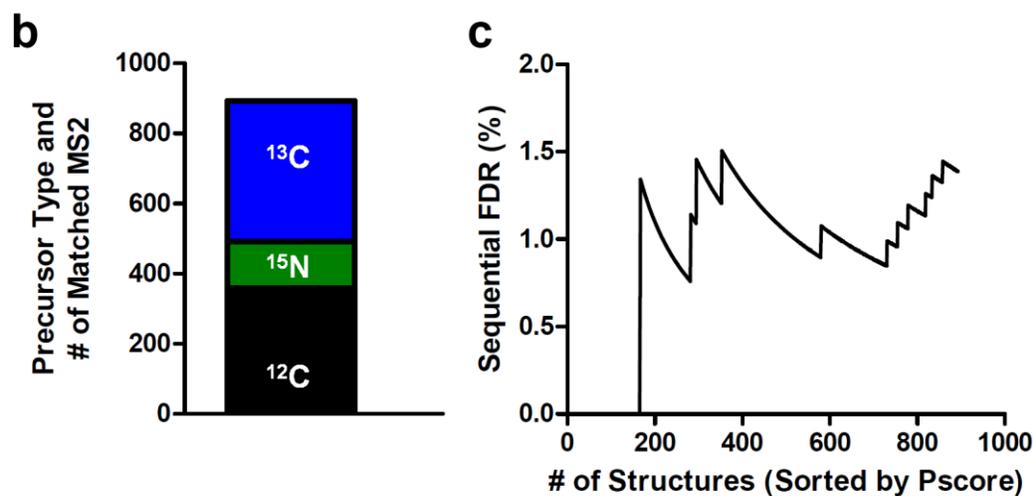
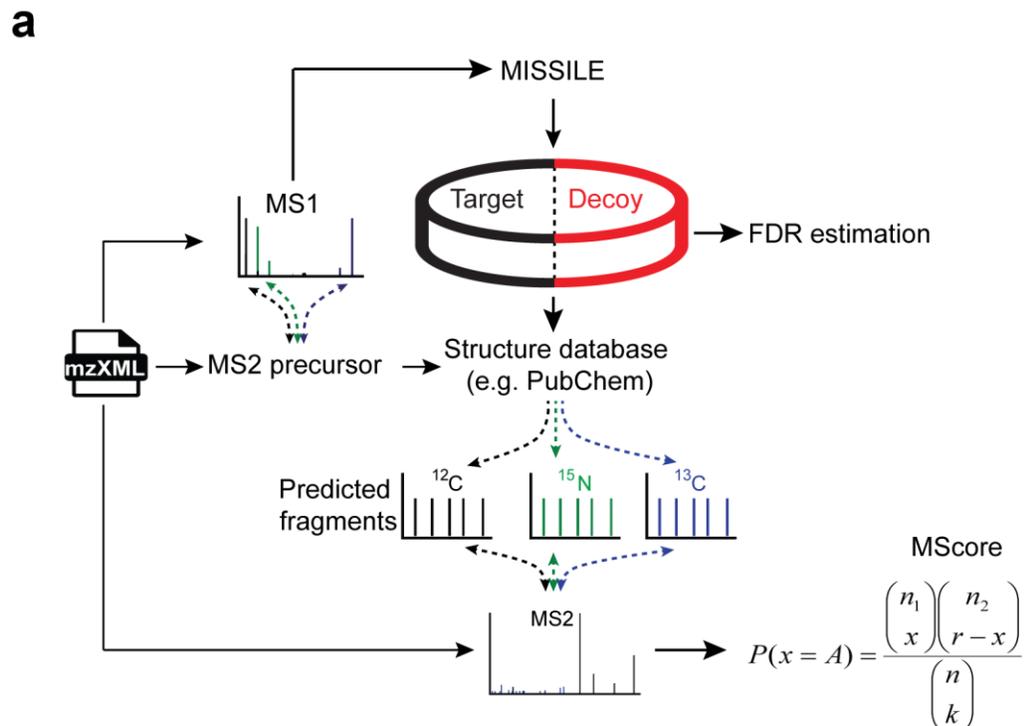
Supplementary Fig. 8. Systematic parameter optimization. (a-b) The number of target MISSILES (detected isotope labeled pairs) that pass both mass-formula and structure databases (e.g., HMDB) increases as the correlation coefficient and relative intensity decrease, respectively. (c-d) The number also increases as mass defect tolerance for pairing and database search tolerance increase. The structure database hits reach a plateau at 0.5 for both peak correlation coefficient and relative intensity, and at 4 ppm for both mass defect tolerance and database search tolerance, respectively. Percent numbers above the data points represent the apparent false discovery rate under those search parameters. These optimized parameters may be column and instrument dependent and should be reassessed for each user's platform.



162
 163 **Supplementary Fig. 9. Number of molecular formulas containing 16 selected elements in**
 164 **the HMDB and YMDB.** We aimed to determine the relative representation of 16 common
 165 biological elements among the HMDB and YMDB. Outside of carbon, hydrogen, nitrogen,
 166 oxygen, phosphorus, and sulfur, the remaining elements were uncommon. The YMDB is
 167 relatively enriched for metabolites containing phosphorus compared to the HMDB. However,
 168 metabolites composed of C, H, N, and O account for the bulk of formula entries.
 169



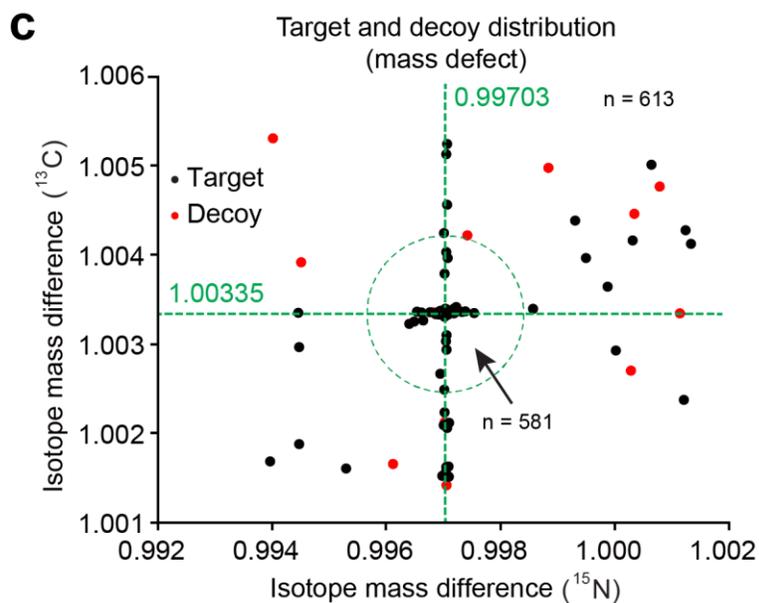
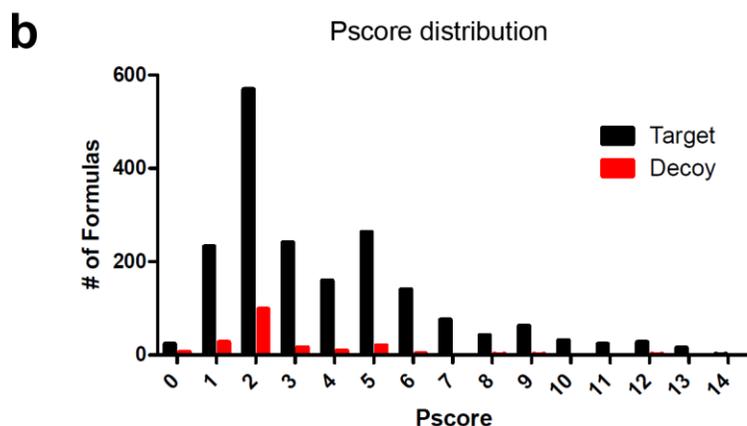
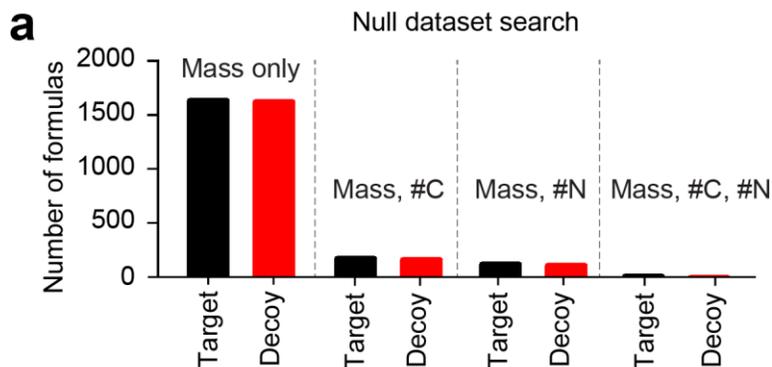
170
 171 **Supplementary Fig. 10. Comparison of the number of formulas between the PubChem,**
 172 **HMDB, and YMDB databases.** We downloaded local copies of these publically available
 173 compound/metabolite databases to determine the degree of redundancy between data sources.
 174 Overall most formulas were contained within Pubchem.
 175
 176



177
178

179 **Supplementary Fig. 11. Scoring of metabolite-spectrum matches.** (a) Putative MISSILE
180 formulas are searched against the target and decoy databases to estimate the rate of false
181 discovery. For each formula JUMPm identifies all known structures sharing the specified
182 chemical formula. For each candidate structure, the tandem mass spectrometry fragments are
183 predicted for the unlabeled and theoretical labeled versions using the CDK implementation of
184 MetFrag. Empirical tandem scans are matched and scored against the predicted fragments
185 using the formula shown (MScore). (b) Proportion of MISSILEs with MS2 scans from each type
186 of isotope label. The data-dependent Top N acquisition may target ^{12}C , ^{15}N , or ^{13}C parent ions.
187 Therefore, only $\sim 1/3^{\text{d}}$ of tandem scans come from carbon-12 metabolite precursors. There are
188 fewer ^{15}N precursors in part because less than half of metabolites contain nitrogen. (c)

189 Cumulative FDR of structures identified by JUMPm. The identified structures were sorted by
 190 Pscore from high to low (left to right). As lower Pscore hits are included, the number of decoys
 191 accumulates, increasing the FDR for that subset of the results. Each “spike” in the plot is due to
 192 a decoy hit. JUMPm uses a Pscore cutoff to trim the data to achieve the user-specified FDR.
 193

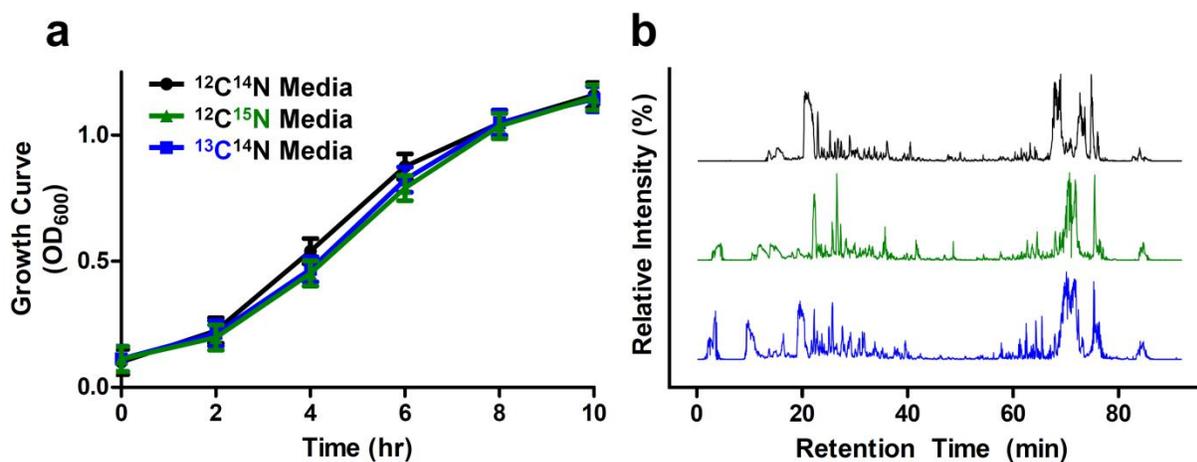


194

195

196 **Supplementary Fig. 12. FDR characterization with labeled yeast.** (a) Effect of formula
197 search criteria on the number of hits and the FDR with a null (negative control) dataset. More
198 stringent criteria (e.g., both mass and carbon number) reduced the number of known spurious
199 hits (target and decoy) without affecting the apparent FDR. (b) Histogram of formulas as a
200 function of Pscore. The number of decoys significantly decreases as Pscore increases. JUMPm
201 can trim the data using the Pscore to achieve the user-specified FDR (e.g., 1%), real dataset
202 used for analysis. (c) Two-dimensional distribution of target and decoy matches by isotope
203 mass differences (581 out of 613 within the green circle), which are derived from the accurate
204 delta mass between labeled and unlabeled ions divided by the integer delta mass (**Online**
205 **Methods**). Decoy matches are randomly scattered in the plot while targets are tightly clustered
206 around the expected mass differences of 1.00335 (carbon) and 0.99703 (nitrogen), real dataset
207 used.

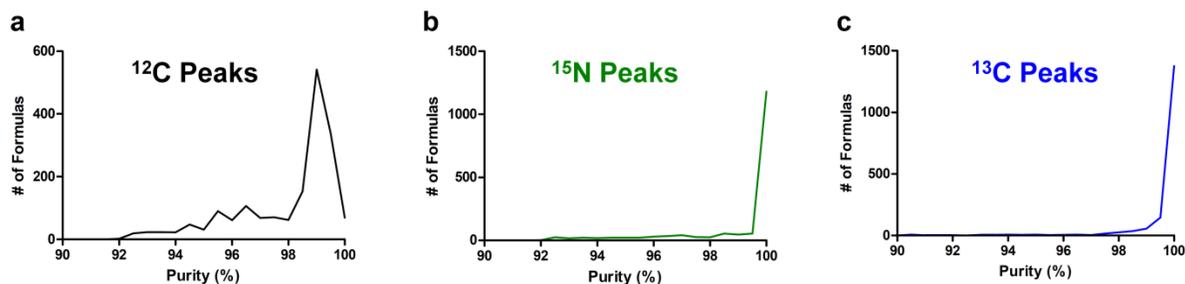
208
209
210
211
212
213
214
215
216
217



218
219
220
221
222
223
224
225
226
227

Supplementary Fig. 13. Stable isotope labeling in yeast. (a) Growth curves for each of the isotopic yeast cultures (unlabeled ¹²C), ¹⁵N, ¹³C). (b) Base-peak LC-MS chromatograms for each of the yeast culture metabolite extracts, C18 column in positive mode (**Online Methods**).

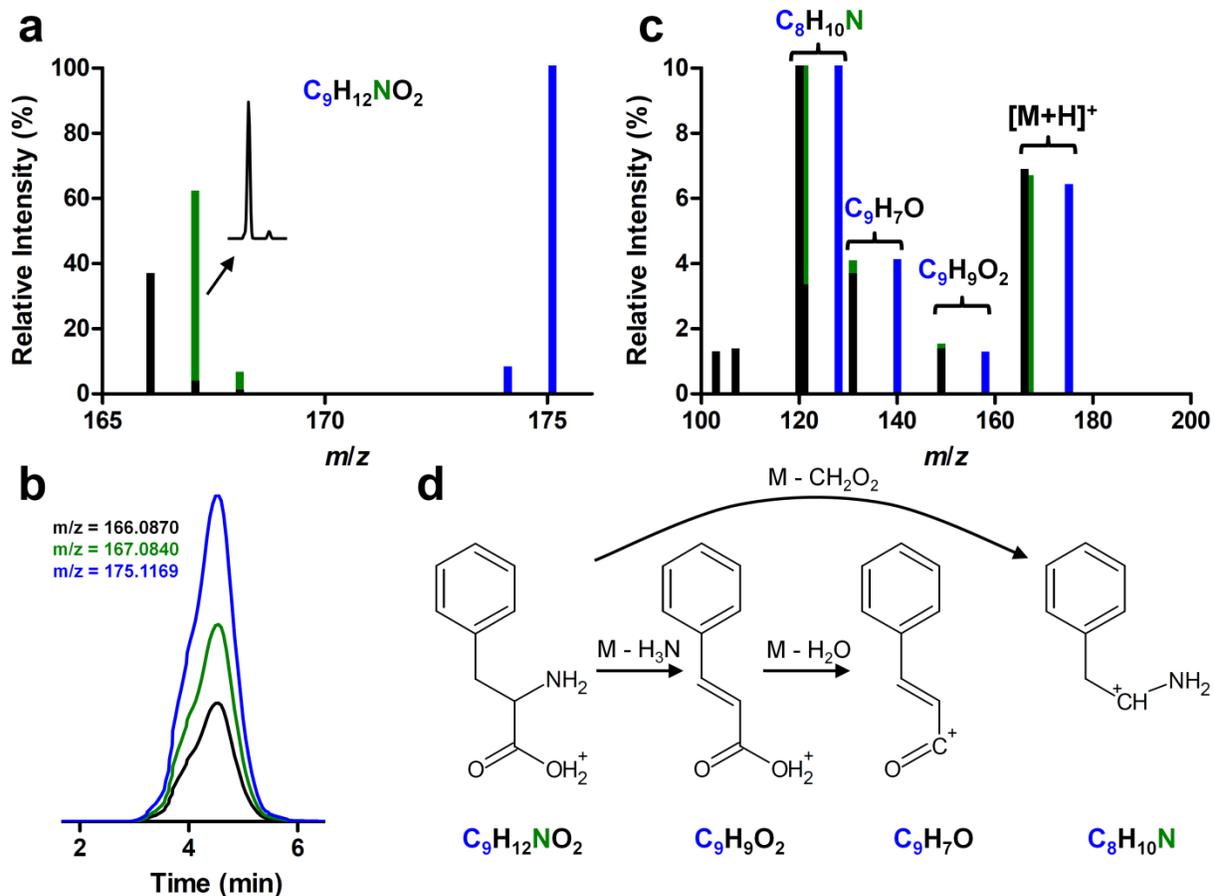
228
229



230
231

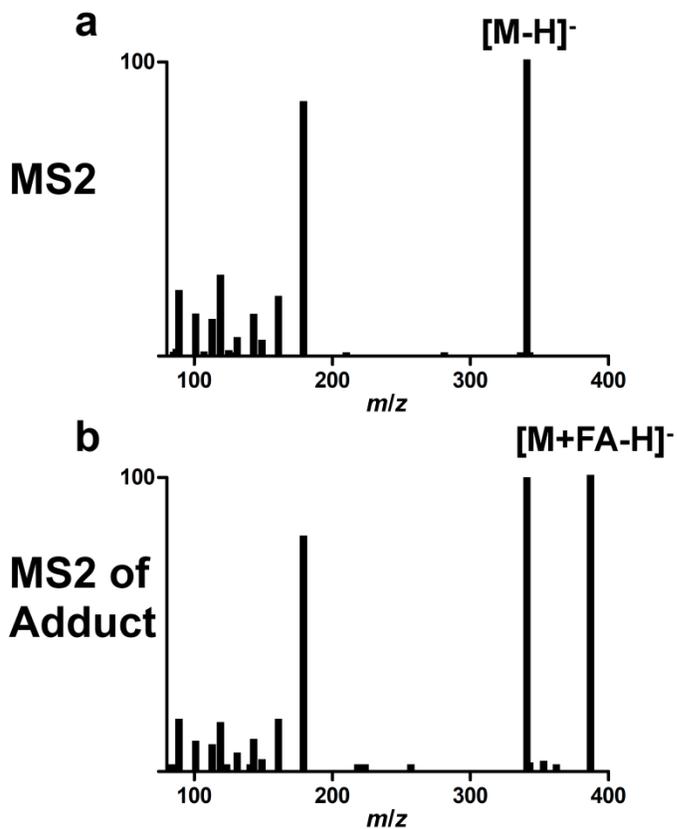
232 **Supplementary Fig. 14. Purity of isotope incorporation in yeast metabolites.** Formulas
233 detected by JUMPm were subjected to isotopic pattern simulation analysis to determine the
234 purity (% labeling) of each labeled ion (**Online Methods**). **(a)** Histogram of unlabeled (^{12}C)
235 peaks (n=1722). The mode purity among unlabeled peaks was 99%, in close agreement with
236 the natural isotopic purity of ^{12}C (98.9%). **(b)** Histogram of ^{15}N labeled peaks (n=1631). The
237 mode purity among ^{15}N labeled peaks was 99%. **(c)** Histogram of ^{13}C labeled peaks (n=1722).
238 The mode purity among ^{13}C labeled peaks was 99%.

239
240
241
242
243
244
245
246
247
248
249
250

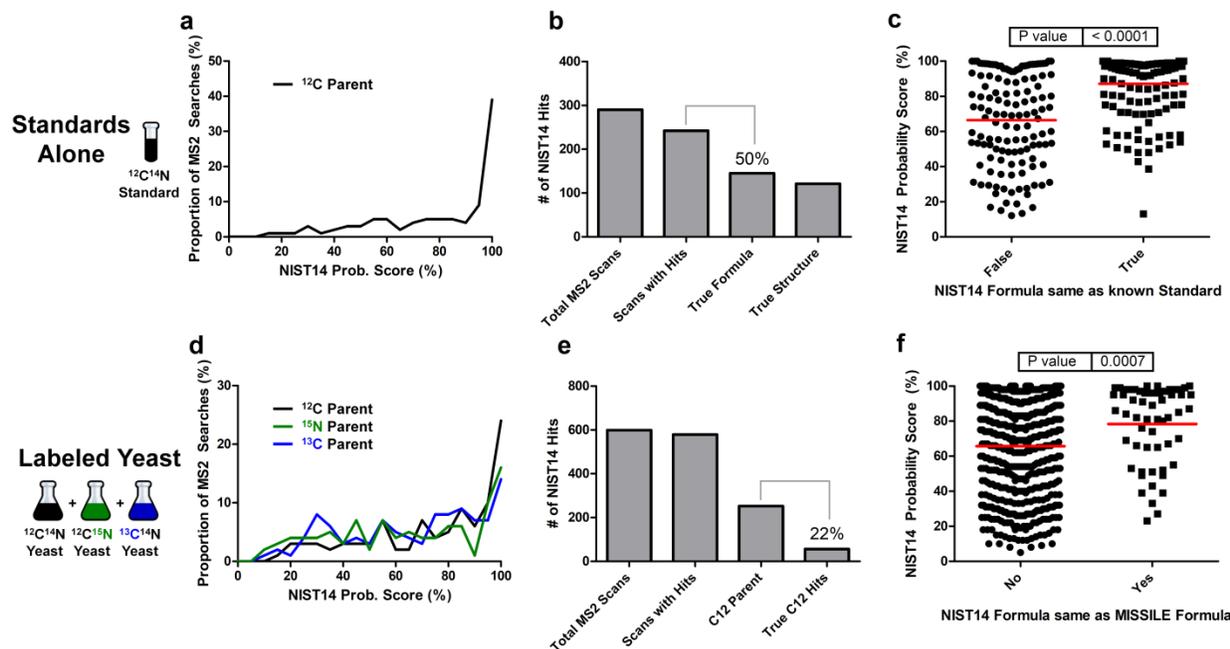


251
 252
 253
 254
 255
 256
 257
 258
 259
 260
 261
 262
 263
 264
 265
 266
 267
 268

Supplementary Fig. 15. MS/MS annotation using stable isotope information. (a) MS1 spectrum of phenylalanine including unlabeled (black), $^{12}C^{15}N$ labeled (green), $^{13}C^{14}N$ labeled (blue) peaks. Arrow points to a profile trace of the raw MS1 data showing complete resolution of the natural isotope A+1 peak and the $^{12}C^{15}N$ labeled peak. (b) Co-eluting labeled ions of phenylalanine in the mixed yeast sample. (c) Overlay (n=3) of separate tandem MS scans for each of the labeled ions in (a). (d) Structural assignment of the phenylalanine fragmentation pathway.



269
 270 **Supplementary Fig. 16. Analysis of a formic acid adduct of trehalose from yeast.** (a) MS2
 271 spectrum of the unlabeled trehalose. The $[M-H]^-$ molecular ion still visible among the fragments.
 272 Characteristic fragments observed below 200 *m/z*. (b) MS2 spectrum of the unlabeled $[M+FA-H]^-$
 273 adduct of trehalose. Parent ion and loss of formic acid molecular ion still visible among the
 274 fragments.
 275
 276
 277
 278



279
 280
 281 **Supplementary Fig. 17. Spectral library search of standards and yeast MS2 scans.** (a-c)
 282 Spectral library search results for the synthetic standard library. (a) Distribution of NIST14
 283 probability scores for MS2 spectra from known standards (**Supplementary Table 5**). (b)
 284 NIST14 MS/MS database results by status. About 50% of spectra searches gave the same
 285 formula as the known standard. Overall, 42% were also the same structure as the standard
 286 compound. (c) NIST14 probability scores among True and False hits for the standard library.
 287 (d-f) Spectral library search results for selected yeast metabolites (n=599). (d) Distribution of
 288 NIST14 probability score for MS2 spectra from unlabeled (^{12}C), ^{15}N , and ^{13}C parent ions. (e)
 289 NIST14 MS/MS database results by status. ^{12}C parents are MS2 scans from unlabeled yeast
 290 precursor ions. “True” hits are those where NIST14 returned the best hit with the correct
 291 chemical formula while “False” hits are those where NIST14 returned the best hit with an
 292 incorrect chemical formula. (f) NIST14 probability scores among True and False hits for the
 293 yeast metabolites.

294
 295
 296
 297
 298
 299
 300
 301
 302
 303
 304
 305
 306