

Supporting Information

Improved annotation of untargeted metabolomics data through buffer modifications that shift adduct mass and intensity

Wenyun Lu, Xi Xing, Lin Wang, Li Chen, Sisi Zhang, Melanie R. McReynolds, Joshua D. Rabinowitz*

Lewis Sigler Institute for Integrative Genomics and Department of Chemistry, Princeton University, Princeton, New Jersey 08544, United States

Correspondence: joshr@Princeton.edu

Table of Content

I. Supporting Methods

1. Chemicals and reagents
2. *S. cerevisiae* cell culture and metabolite extraction
3. Metabolite extraction from mouse liver
4. LC-MS
5. Data analysis

II. Supporting Figures and Tables

Figure S1. Examples showing the statistical distributions of the relative $\Delta m/z$ errors (in ppm) for all the pairwise peaks matching the expected mass differences of (A) ^{13}C - ^{12}C , (B) $\text{C}_2\text{H}_4\text{O}_2$, and (C) Na-H.

Figure S2. The intensity ratio $I_{\text{adduct}}/I_{[\text{M}+\text{H}]^+}$ for selected known adducts in positive ion mode in liver extract.

Figure S3. The intensity ratio of $I_{\text{adduct}}/I_{[\text{M}-\text{H}]^-}$ for selected known adducts in negative ion mode in liver extract.

Figure S4. Example of UDP-D-glucose showing that some metabolite ions behave like fragments when examined by in-source CID due to the presence of abundant adduct ions.

Figure S5. Comparison of the annotation of “Metabolite” for *S. cerevisiae* using PAVE and BMW.

Figure S6. MS2 and labeling analysis of the two peaks in the chromatogram of aspartate.

Figure S7. Annotation problem results from two ions with similar masses that are not fully resolved.

Figure S8. Total ion chromatogram (TIC) of liver extract and extracted ion chromatograms (EIC) of selected metabolites in Buffer-1 and Buffer-2.

Table S1. Parameter setting for “Isotope” and “Adduct” annotation: $\Delta m/z$ and tolerance, RT tolerance, lower bound and upper bound for $I_{\text{isotope}}/I_{\text{metabolite}}$ and $I_{\text{adduct}}/I_{\text{metabolite}}$ (ac0c00985_si_002.xlsx).

Table S2. BMW workflow for *S. cerevisiae* data analysis, showing the datasets used for different annotation steps (Word table).

Table S3. PAVE workflow for *S. cerevisiae* data analysis, showing the datasets used for different annotation steps (Word table).

Table S4. Results of *S. cerevisiae* data annotation using PAVE and BMW (ac0c00985_si_003.xlsx).

Table S5. Validation of BMW using known metabolite, adduct and fragment peaks in liver data (ac0c00985_si_004.xlsx).

Table S6. Results of liver data annotation (ac0c00985_si_005.xlsx).

Table S7. Summary of liver LC-MS peak annotation (Word table).

Supporting Methods

1. Chemicals and reagents

HPLC-grade water (W6), acetonitrile (A955), and methanol (A456) were from Thermo Fisher. Other components for the LC mobile phase are unlabeled ammonium hydroxide (A669S-500, 28.0 to 30.0 w/w %, Fisher), unlabeled ammonium acetate (238074, $\geq 97\%$, Sigma), unlabeled ammonium formate (516961, $\geq 99.995\%$, Sigma), ^{15}N -ammonium hydroxide solution (488011, $\sim 3\text{N}$ in H_2O , 98% atom ^{15}N , Sigma), ^{15}N -ammonium acetate (363006, 98 atom % ^{15}N , Sigma). U- ^{13}C -Glucose (CLM-1396, 99%) and $(^{15}\text{NH}_4)_2\text{SO}_4$ (NLM-713, 99%) were obtained from Cambridge Isotope Laboratories. All other chemicals were obtained from Sigma.

Metabolite standards were obtained from Sigma (St. Louis, MO), Avanti (Alabaster, AL), or MetaSci (Toronto, Canada).

2. *S. cerevisiae* cell culture and metabolite extraction

The *S. cerevisiae* growth medium contains the following: YNB without vitamin and amino acids (Sunrise, #1524) 1.7 g/L, biotin 0.002 mg/L, glucose 20 g/L, $(\text{NH}_4)_2\text{SO}_4$ 5 g/L. Glucose and $(\text{NH}_4)_2\text{SO}_4$ are either unlabeled, or ^{13}C - or ^{15}N -labeled resulting in four conditions: unlabeled, ^{13}C , ^{15}N and $^{13}\text{C}/^{15}\text{N}$.

The YNB without vitamin and amino acids (Sunrise, #1524) contains the following components according to the company's website:

Boric acid (0.5 mg/L)

Calcium chloride dihydrate (100 mg/L)

Copper (II) sulfate pentahydrate (0.04 mg/L)

Iron (III) chloride (0.2 mg/L)

Magnesium sulfate anhydrous (500 mg/L)

Manganese sulfate monohydrate (0.4 mg/L)

Potassium iodide (0.1 mg/L)

Potassium phosphate monobasic anhydrous (1000 mg/L)

Sodium chloride (100 mg/L)

Sodium molybdate (0.2 mg/L)

Zinc sulfate monohydrate (0.4 mg/L)

S. cerevisiae strain FY4 was grown at 30 °C in growth medium as outlined above. Labeling was carried out for > 10 generations. To obtain a procedure blank sample for unlabeled culture, cells were harvested at OD₆₀₀ of 0.02 by filtering 10 mL culture onto a 50 mm nylon membrane filter (0.45 µm pore size, Millipore), which was immediately transferred into -20°C extraction solvent (1 mL 40:40:20:0.5% FA, acetonitrile/methanol/water/formic acid) in a Petri dish. The dish was kept at -20°C for 5 min. Then 84 µL 15% NH₄HCO₃ (w:v) was added to neutralize the samples. The final solution was kept at -20 °C for 15 min and the resulting mixture was transferred into an Eppendorf tube and spun down at 16,000 g for 15 min at 4°C. The supernatant was taken as the final extract. In addition, for all cultures, they were allowed to continuously grow until OD₆₀₀ = 0.80 and 10 ml cultures were then harvested and metabolite extracted using same volume of extraction solvent as above. In total, five extracts were generated and analyzed via LC-MS in addition to the solvent blank samples: unlabeled extract at OD₆₀₀=0.02, unlabeled extract at OD₆₀₀=0.80, ¹⁵N extract at OD₆₀₀=0.80, ¹³C extract at OD₆₀₀=0.80, and ¹³C/¹⁵N extract at OD₆₀₀=0.80.

3. Metabolite extraction from mouse liver

Twelve-month-old female wild-type C57BL/6 mice (The Jackson Laboratory, Bar Harbor, ME) on normal diet were sacrificed by cervical dislocation and tissues quickly dissected and snap frozen in liquid nitrogen with precooled Wollenberger clamp. This avoids further metabolite turnover and helps to protect unstable compounds. Frozen samples from liquid nitrogen were then transferred to -80°C freezer for storage. To extract metabolites, frozen liver tissue samples were first weighed (~ 20 mg each) and transferred to 2.0 mL round-bottom Eppendorf Safe-Lock tubes on dry ice. Samples were then ground into powder using 6.5mm YSZ balls (4039GM-S065, Inframat® Advanced Materials LLC, Manchester, CT) with a cryomill machine (Retsch, Newtown, PA) for 30 seconds at 25 Hz, and maintained at cold temperature using liquid nitrogen. For every 25 mg tissues, 1 mL 40:40:20 acetonitrile:methanol:water with 0.5% formic acid was added to the tube, vortexed for 10 seconds, and allowed to sit on ice for 10 minutes. Then 84 µL 15% NH₄HCO₃ (w:v) was added and vortexed to neutralize the samples. The samples were allowed to sit on ice for another 20 minutes and then centrifuged at 16,000 g for 25 min at 4°C. The supernatants were transferred to another Eppendorf tube and centrifuged at 16,000 g for another 25 min at 4°C. The supernatants were transferred to glass vials for LC-MS analysis. A

total of two rounds of centrifugation is necessary to remove small tissue particles which will otherwise clog the LC columns. In addition, a procedure blank sample was generated identically without tissue.

In a separate experiment, fresh made liver extract was prepared as above, transferred into a glass vial and loaded into the autosampler maintained at 5 °C, and continuously analyzed over a 24-hour period to evaluate metabolite stability, as well as discovery of any new chemicals produced during this time range.

4. LC-MS

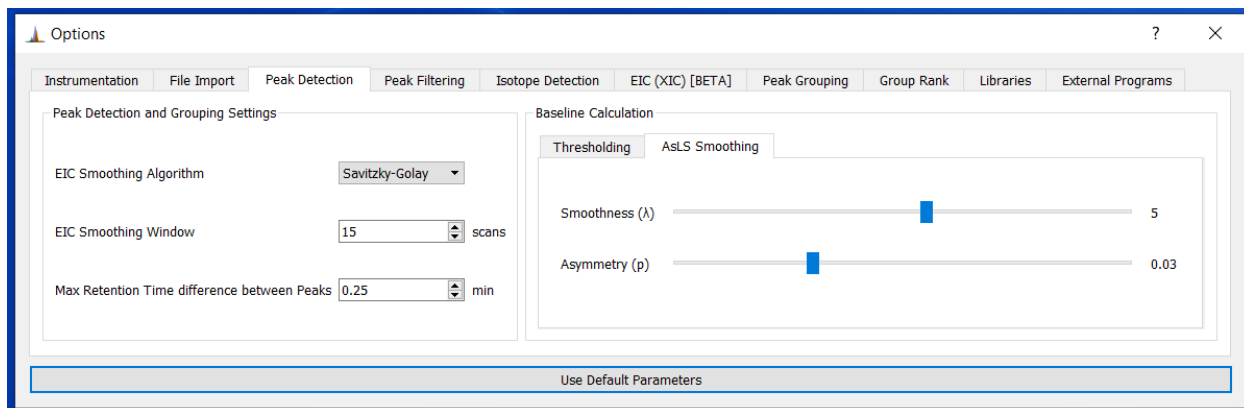
LC separation was achieved using a Vanquish UHPLC system (Thermo Fisher Scientific) with an Xbridge BEH Amide column (150×2mm, 2.5 µm particle size; Waters, Milford, MA). Solvent A is 95:5 water: acetonitrile with 20 mM ammonium acetate and 20 mM ammonium hydroxide at pH 9.4, and solvent B is acetonitrile. This buffer system was designated as “Buffer-1” which contains 40 mM NH₄⁺ and 20 mM CH₃COO⁻. For modified buffer (“Buffer-2”), solvent A is 95:5 water:acetonitrile with 10 mM NH₄OH, 10 mM ¹⁵NH₄OH, 10 mM CH₃COO⁻¹⁵NH₄ and 10 mM HCOONH₄ resulting in 20 mM ¹⁴NH₄⁺, 20 mM ¹⁵NH₄⁺, 10 mM CH₃COO⁻ and 10 mM HCOO⁻. Solvent B is acetonitrile. The gradient is 0 min, 90% B; 2 min, 90% B; 3 min, 75%; 7 min, 75% B; 8 min, 70%, 9 min, 70% B; 10 min, 50% B; 12 min, 50% B; 13 min, 25% B; 14 min, 25% B; 16 min, 0% B, 20.5 min, 0% B; 21 min, 90% B; 25 min, 90% B. Total running time is 25 min at a flow rate of 150 µl/min. For all experiments, 5 µl of extract was injected with column temperature set to 25 °C.

The extracts were analyzed in separate runs of either positive ion mode or negative ion mode (e.g. no polarity switch) on a Q-Exactive Plus mass spectrometer. All samples were analyzed with a MS1 scan range of *m/z* 70-1000 with the highest resolving power setting (160,000 at *m/z* 200). Samples were analyzed with either in-source CID turned off (0 eV), or with in-source CID at 5 or 10 eV. It is necessary to set the user role as “Advanced” in order to access the in-source CID feature. Other MS parameters are as follows: sheath gas flow rate, 28 (arbitrary units); aux gas flow rate, 10 (arbitrary units); sweep gas flow rate, 1 (arbitrary units); spray voltage, 3.3 kV; capillary temperature, 320°C; S-lens RF level, 65; AGC target, 3E6 and maximum injection time, 500 ms. Samples were ran in triplicate.

When needed, targeted MS2 scans were performed using the PRM function with the parent ions in the inclusion list. MS2 spectra were collected at HCD energy of 25 eV with other instrument setting being, resolution 17500, AGC target 1e6, Maximum IT 500 ms, isolation window 1.5 m/z.

5. Data analysis

Thermo LC-HRMS raw data were converted to mzXML format using the “msconvert” tool from ProteoWizard (<http://proteowizard.sourceforge.net/tools.shtml>). To obtain a peak list of interest from the *S. cerevisiae* data at OD₆₀₀=0.80, or from the liver data, we used the El-Maven software package (<https://resources.elucidata.io/elmaven>). The parameters for peak picking in El-Maven were: mass domain resolution 10 ppm, time domain resolution 50 scans, minimum intensity range 1000, minimum peak width 10 scans, minimum signal/baseline ratio 3, minimum good peaks 3 (i.e. the peak should appear in all three replicates), EIC smoothing window 15 scan. An abundance threshold of peak height > 10³ ion counts, and a minimum peak width of 10 scans were applied to filter out low intensity and/or noisy peaks. A setting of minimum signal/baseline ratio of 3 was applied to remove peaks with high baseline. Screenshots of the peak picking parameters are shown below. Duplicated peaks were then removed by using *m/z* and RT tolerances of 10 ppm and ±0.1 min, respectively.



Peak Detection ? X

Feature Detection Selection Group Filtering Method Summary

Automated Feature Detection: Find peaks by slicing m/z and retention time space

Mass Domain Resolution: 10.000 ppm Limit m/z Range: 0.00 - 1000000000.00

Time Domain Resolution: 50 scans Limit Time Range: 0.00 - 1000000000.00

Auto Detect And Ignore Isotopes Limit Intensity Range: 1000.00 - 9999999999.00

Compound Database Search: Limit slices to set of known m/z and retention time values

Select Database: KNOWNS

EIC Extraction Window (+/-): 10.00 ppm

Match Retention Time (+/-): 1.00 min

Limit Number of Reported Groups per Compound: 10 best

Match Fragmentation

Fragment Mass Tolerance: 20.00 ppm

Match at least X peaks: 3

Scoring Algorithm: Hypergeometric Score

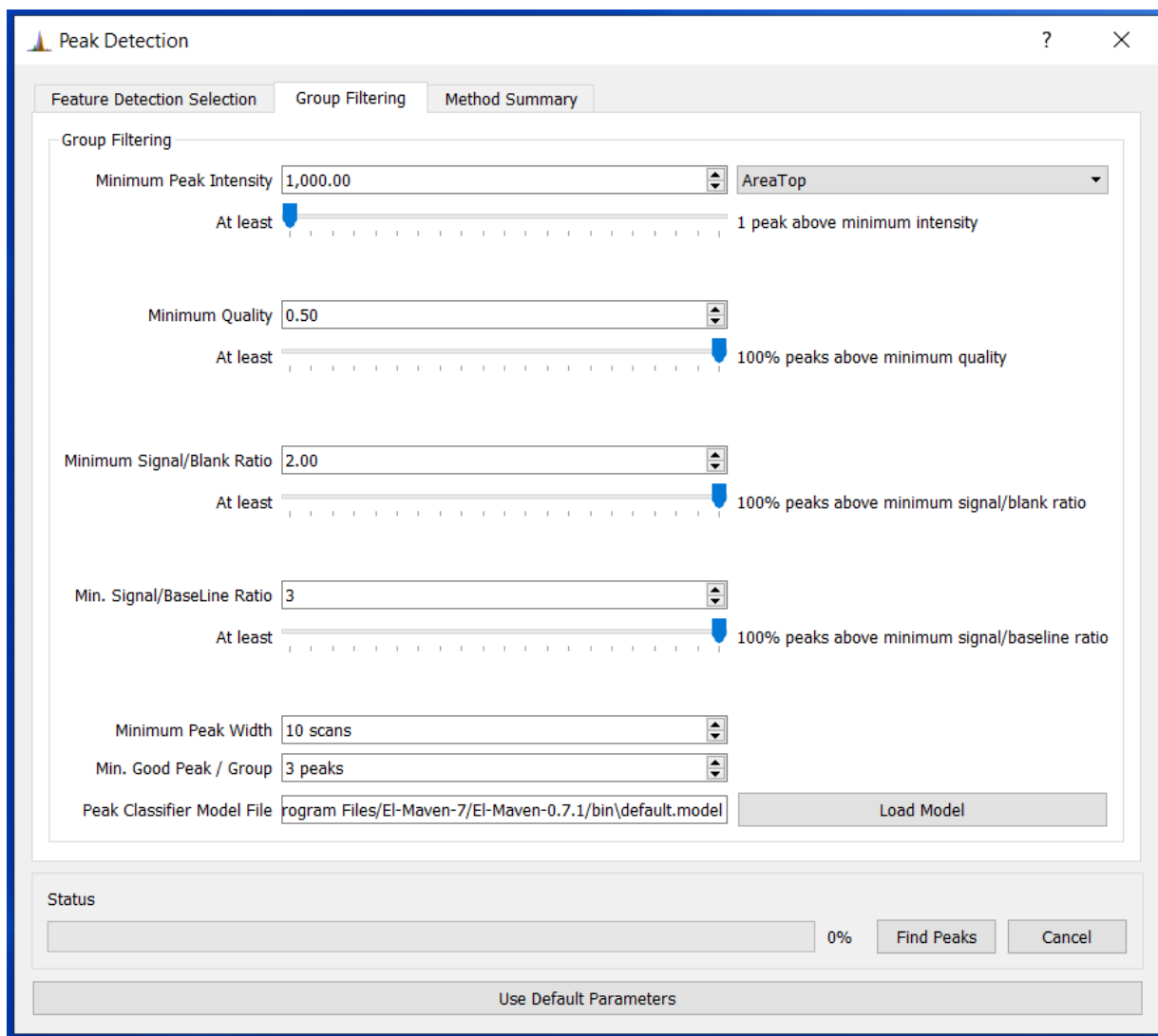
Minimum Score: 0.00

Report Isotopic Peaks

Isotope Detection Options

Status: 0% Find Peaks Cancel

Use Default Parameters



The BMW workflow was implemented in MATLAB. From the peak list table with information on m/z , RT, and signal intensity, different peaks were then annotated with different “modules”, using a set of rules defined in the main text and summarized in Table 1. This often involves numerical calculations as well as matrix transformation, which are handled well by MATLAB. In some cases, it is necessary to read additional signal intensity such as that of $[M+^{15}\text{NH}_4]^+$ (which is not in the initial peak list) in Buffer-2 raw data in order to accurately annotate NH_4^+ adducts. Note that a single peak may fit multiple annotation categories, e.g. an isotope peak can also be an adduct. The count for each category depends on the order of applying the different rules. This order does not, however, impact the set of peaks that are putative metabolites.

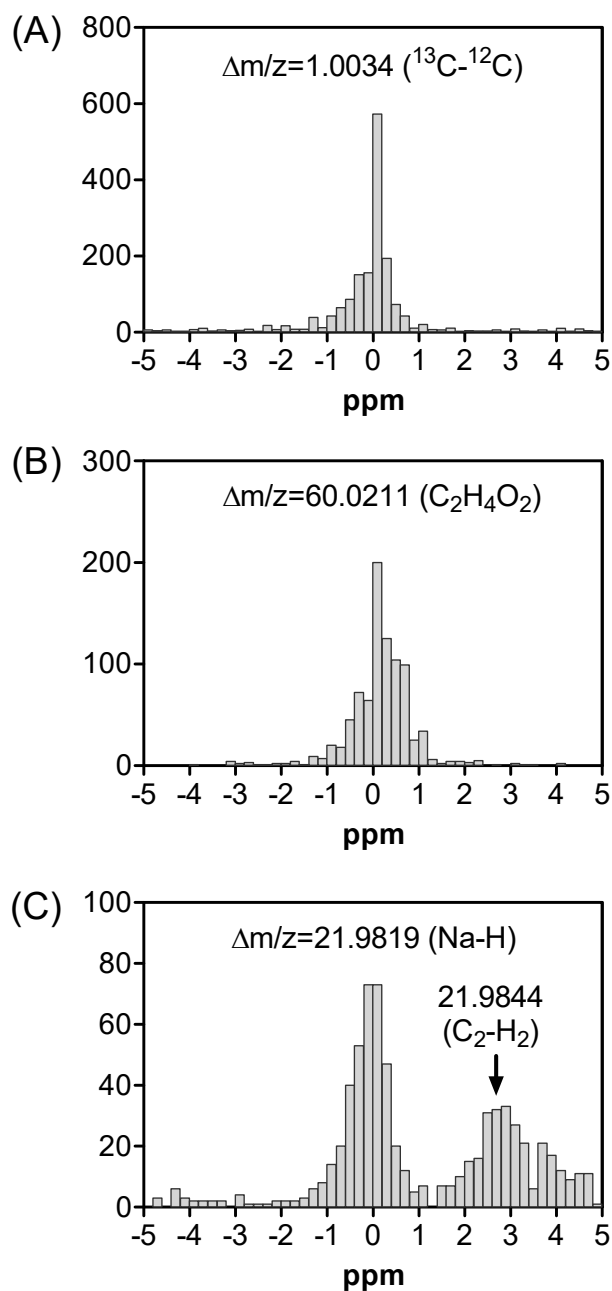


Figure S1. Examples showing the statistical distributions of the relative $\Delta m/z$ errors (in ppm) for all the pairwise peaks matching the expected mass differences of (A) $^{13}\text{C}-^{12}\text{C}$, (B) $\text{C}_2\text{H}_4\text{O}_2$, and (C) Na-H. For this analysis, the retention time differences of the peak pairs are restricted to be within ± 0.1 min. The parameter settings of $\Delta m/z$ (in ppm) cutoff was determined based on these plots. For [Na-H] adduct, it is necessary to use a narrow tolerance (± 1.5 ppm) so as to minimize the contamination from [C_2-H_2] that has very similar $\Delta m/z$ (21.9844 vs. 21.9819). For all other cases, a tolerance of ± 3 ppm was used as it balances well between accuracy and coverage. Y-axis is number of peaks.

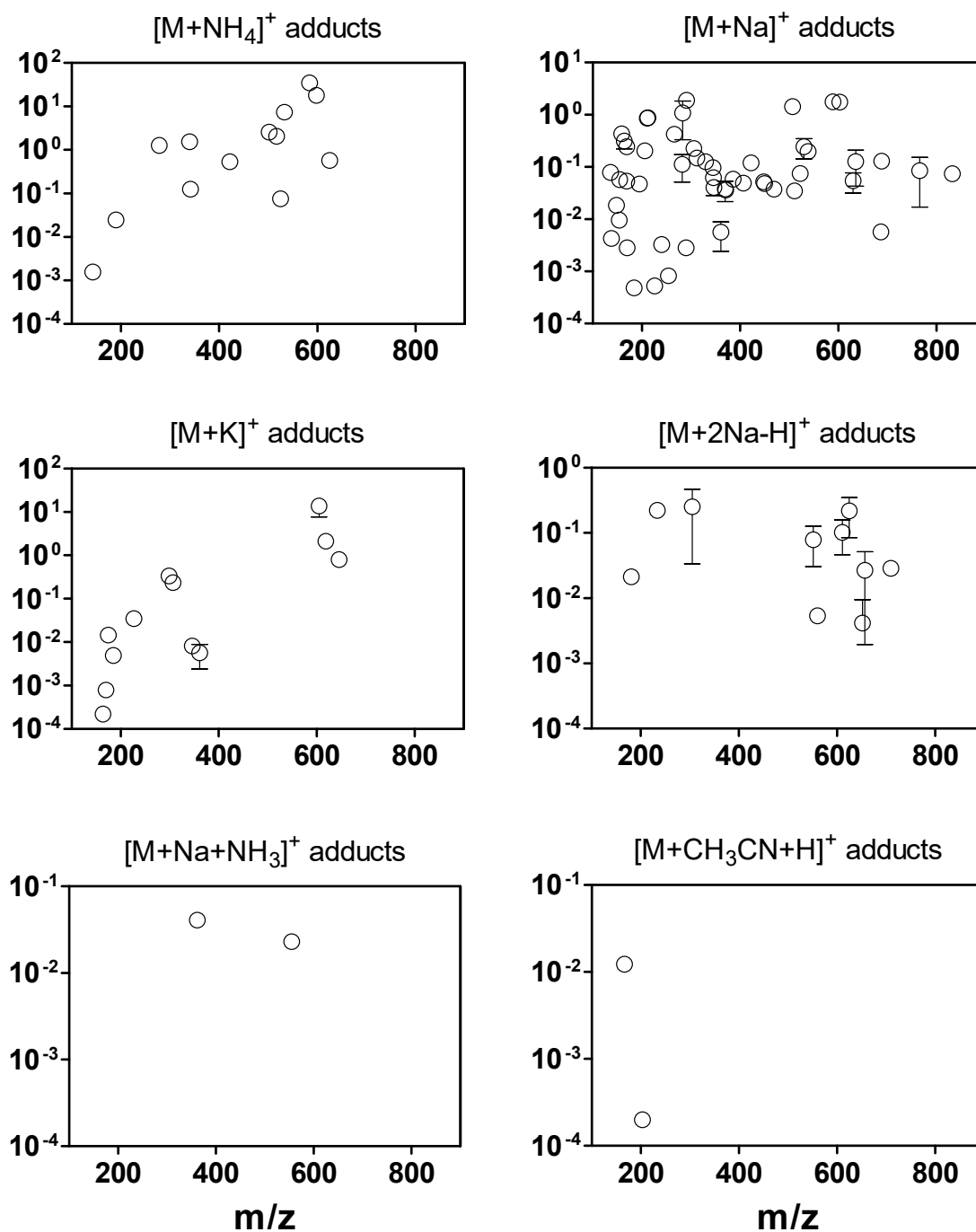


Figure S2. The intensity ratio $I_{\text{adduct}}/I_{[\text{M}+\text{H}]^+}$ for selected known adducts in positive ion mode in the liver extract. This was used to derive the lower bound and upper bound for $I_{\text{adduct}}/I_{[\text{M}+\text{H}]^+}$ ratio as a constraint for adduct annotation in positive ion mode (Table S1).

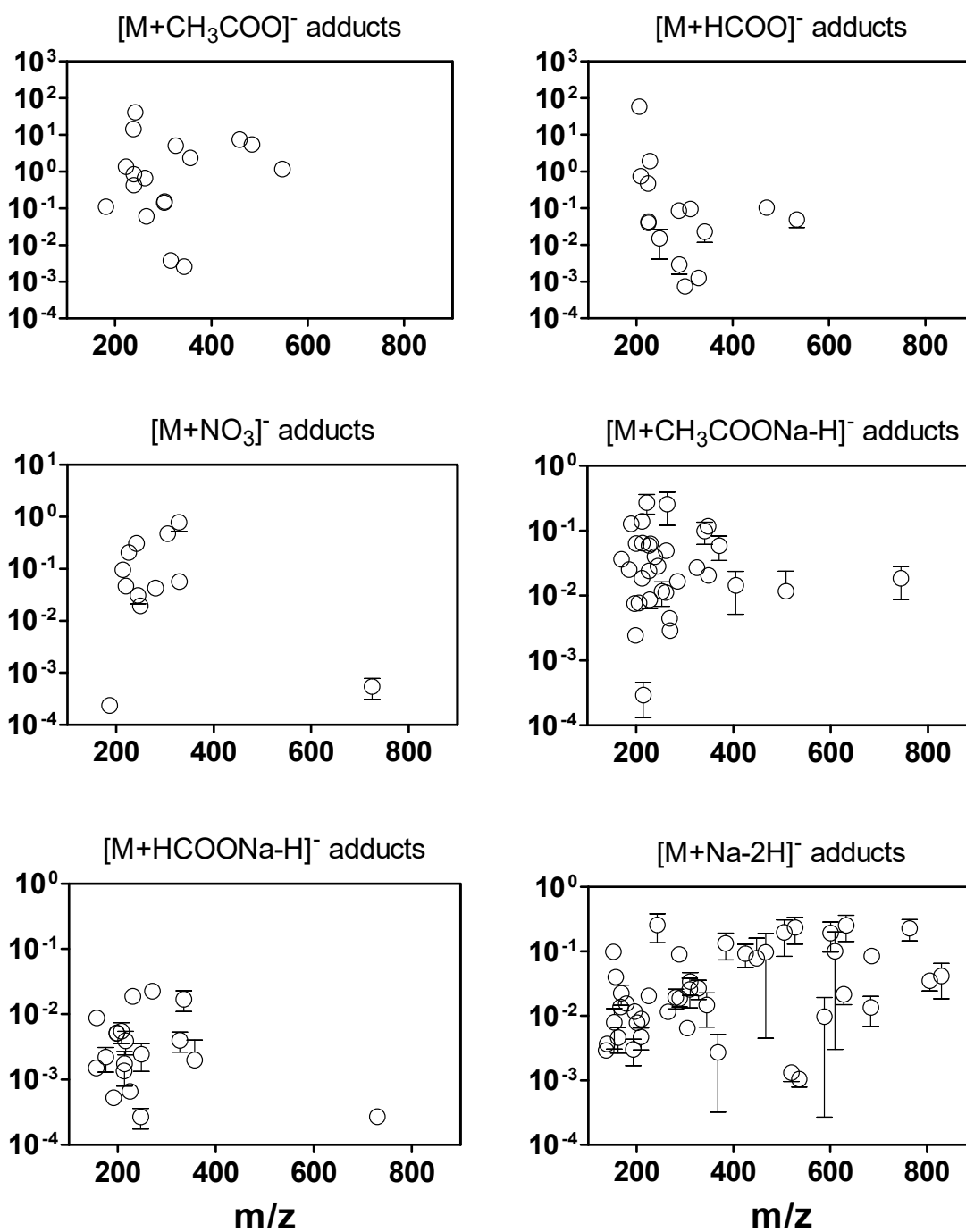


Figure S3. The intensity ratio of $I_{\text{adduct}}/I_{[\text{M-H}]^-}$ for selected known adducts in negative ion mode in the liver extract. This was used to derive the lower bound and upper bound for $I_{\text{adduct}}/I_{[\text{M-H}]^-}$ ratio as a constraint for adduct annotation in negative ion mode (Table S1).

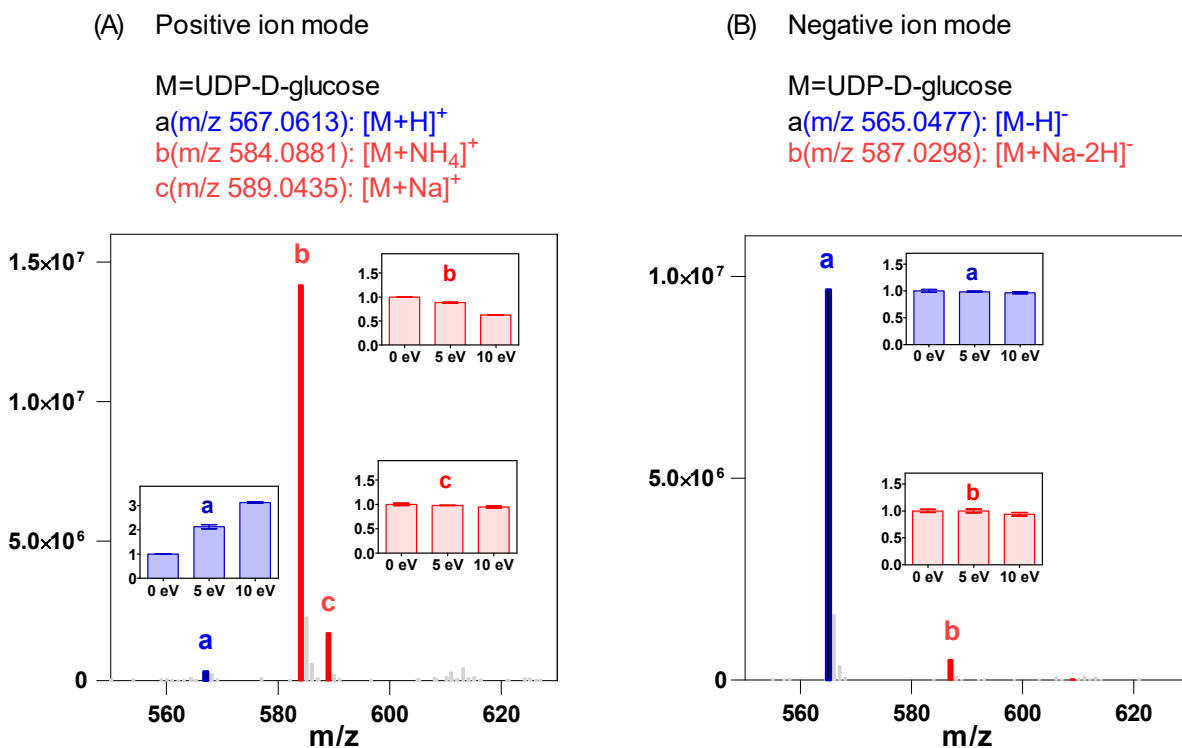


Figure S4. Example of UDP-D-glucose showing that some metabolite ions behave like fragments when examined by in-source CID due to the presence of abundant adduct ions. (A) In positive ion mode, the ion counts for [M+H]⁺ (a), [M+NH₄]⁺ (b), [M+Na]⁺ (c) are ~3e5, ~1.4e7, ~1.6e6, respectively. Applying in-source CID at either 5 or 10 eV increases the intensity of [M+H]⁺, presumably due to the contribution from the fragmentation of [M+NH₄]⁺ or [M+Na]⁺. (B) In negative ion mode, the ion count for [M-H]⁻ (a), [M+Na-2H]⁻ (b) are ~9.6e6, ~5e5, respectively. Applying in-source CID at either 5 or 10 eV decreases the intensity of [M-H]⁻, which is “normal” for metabolite ions. In this case, the preferred ionization mode for UDP-D-glucose is negative mode. The fragment-like behavior is limited to positive mode, where the [M+H]⁺ signal is weak.

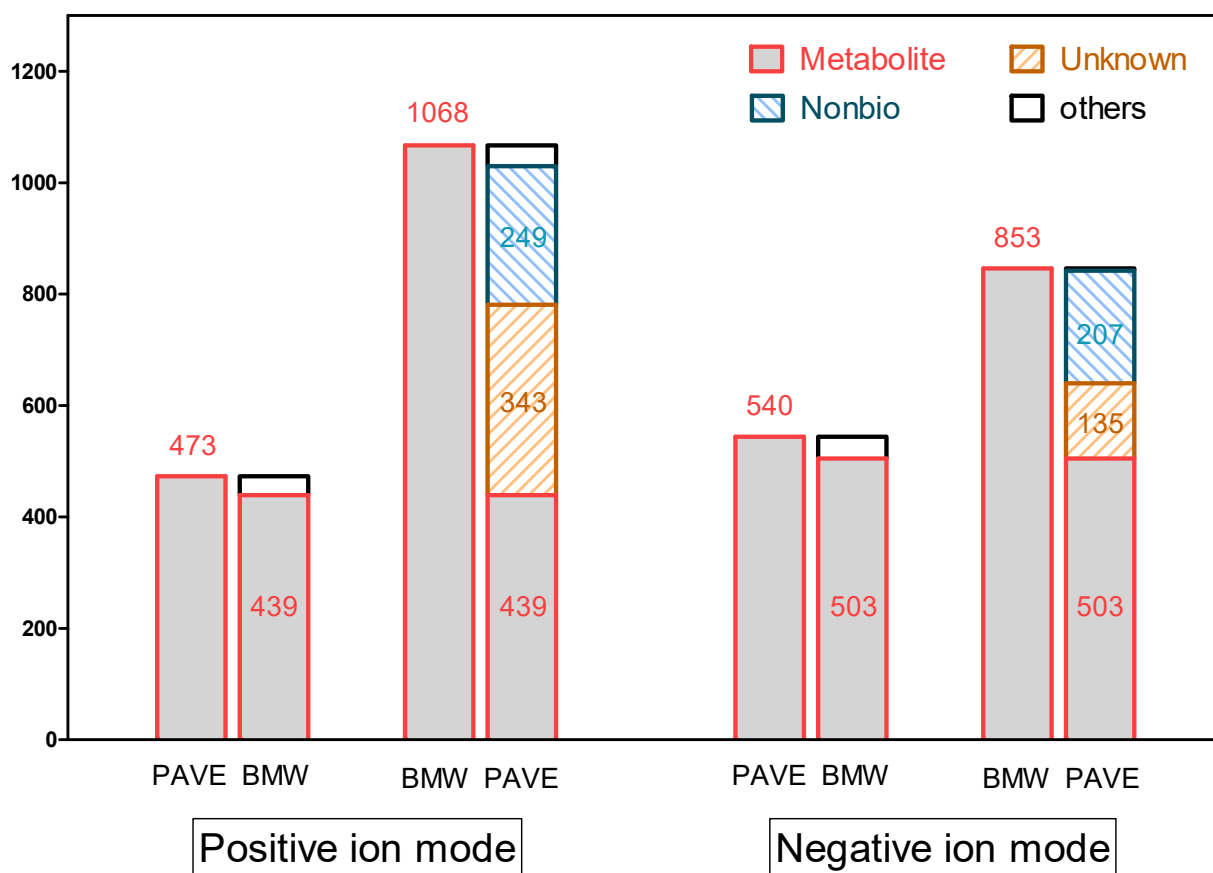


Figure S5. Comparison of the annotation of “Metabolite” for *S. cerevisiae* using PAVE and BMW. Metabolite annotation in PAVE is more stringent than BMW because of the requirement for the molecular formula to match both the metabolite mass and C/N counts. BMW was able to identify 93% of metabolites found in PAVE (439 out of 473 in positive mode, and 503 out of 540 in negative mode). “Unknown” in PAVE refers to those peaks with C/N numbers not matching database formula. This includes those “Reaction product” such as formyl-serine. “Nonbio” in PAVE refers to those peaks without clear labeling pattern so that C/N numbers cannot be determined.

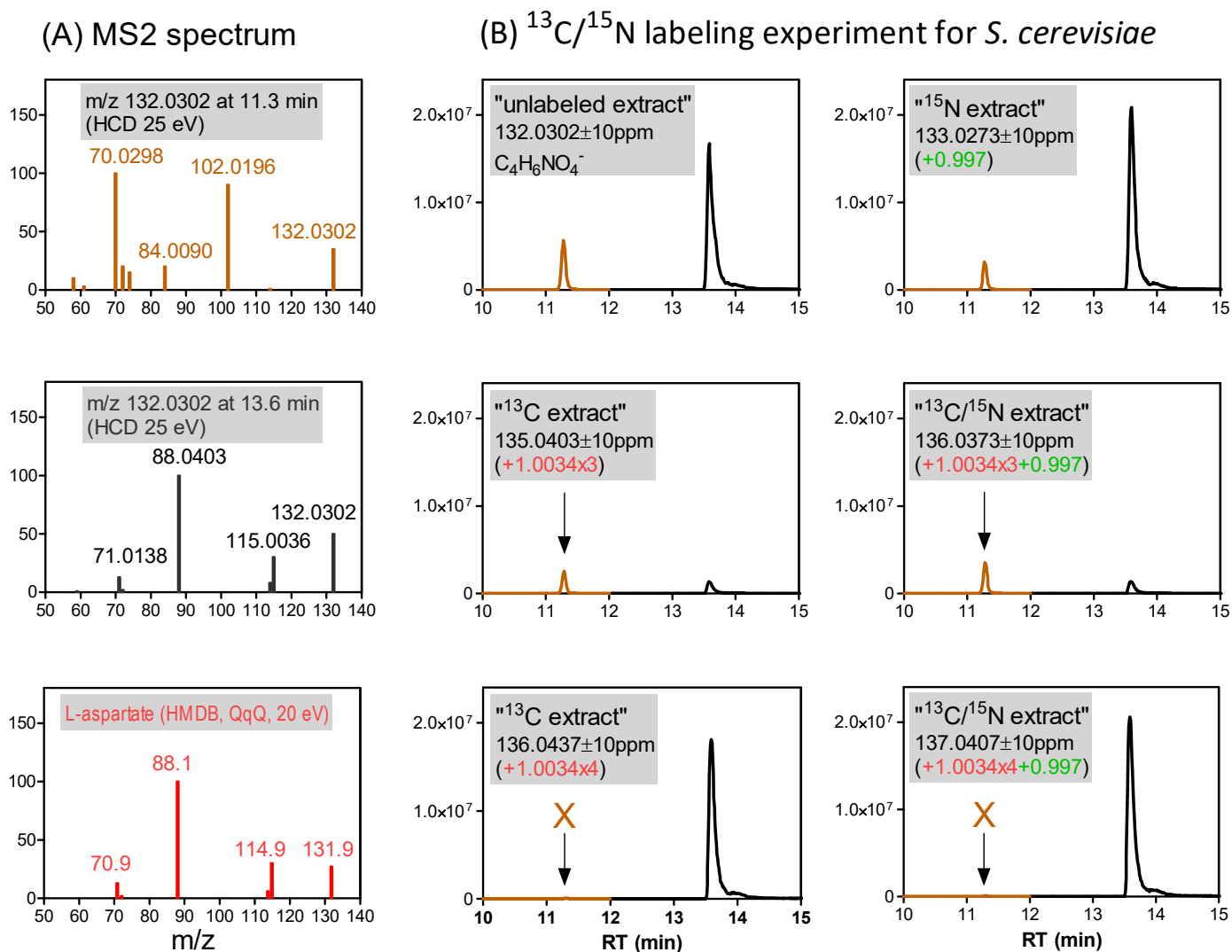


Figure S6. MS2 and labeling analysis of the two peaks in the chromatogram of aspartate.

(A) MS2 spectrum shows that the two peaks correspond to different species and the peak at 13.6 min is confirmed to be aspartate, based on MS2 spectrum matching to HMDB database and RT match with standard. (B) Labeling experiments show that the two species have different labeling patterns. Top panels are at the masses of aspartate (left) and ^{15}N -aspartate (right). Middle panels are at the masses of $^{13}\text{C}_3$ -aspartate (left) and $^{15}\text{N},^{13}\text{C}_3$ -aspartate (right). Lower panels are at the masses of $^{13}\text{C}_4$ -aspartate (left) and $^{15}\text{N},^{13}\text{C}_4$ -aspartate (right). All four carbons can get labeled for aspartate (13.6 min), while only three carbons can get labeled for the peak at 11.3 min (note the absence of the $^{13}\text{C}_4$ labeled forms, marked by "X").

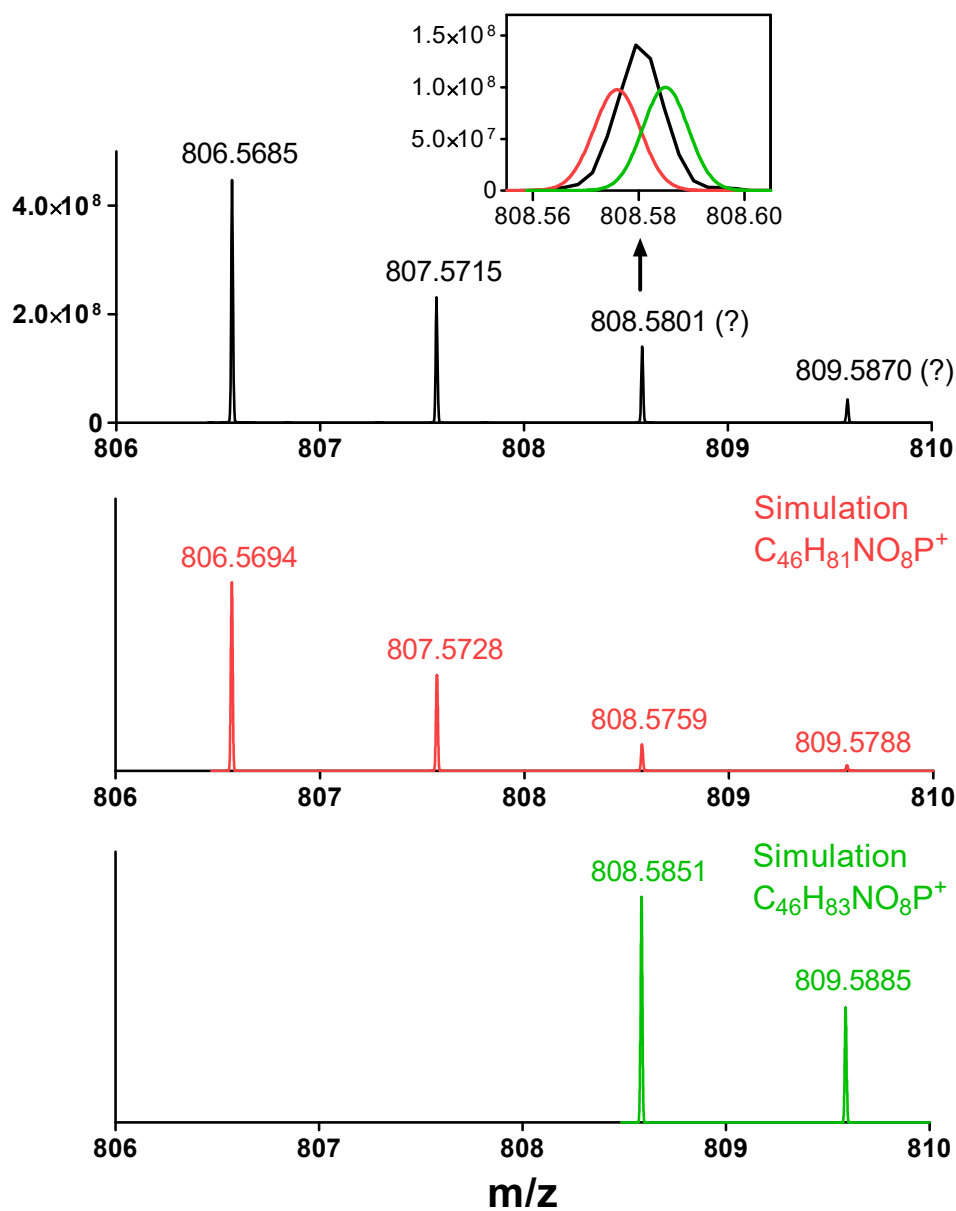
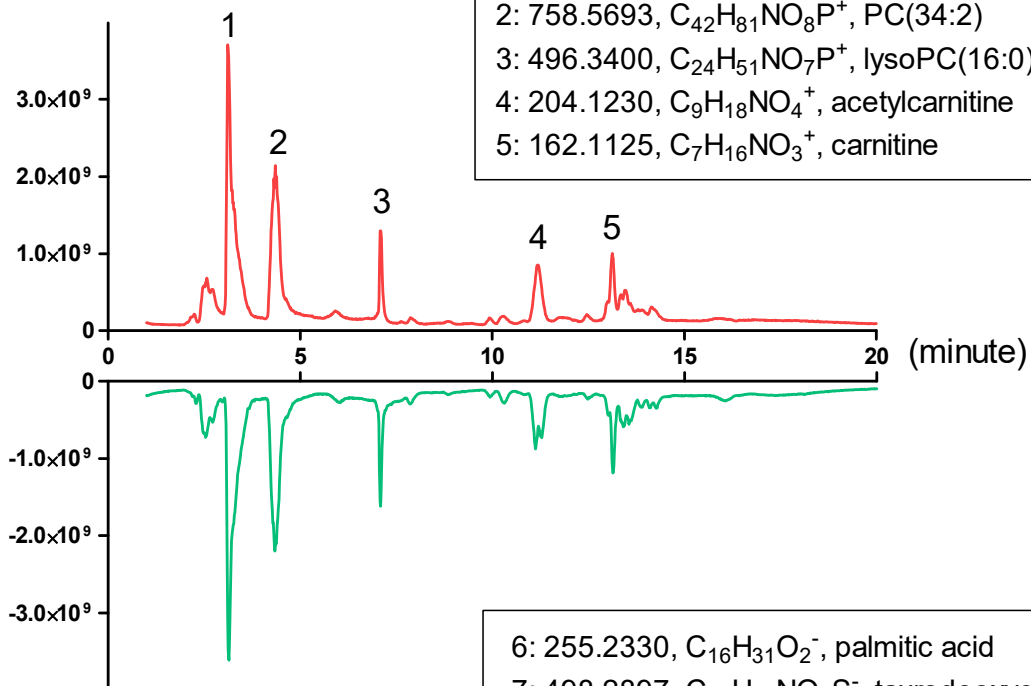
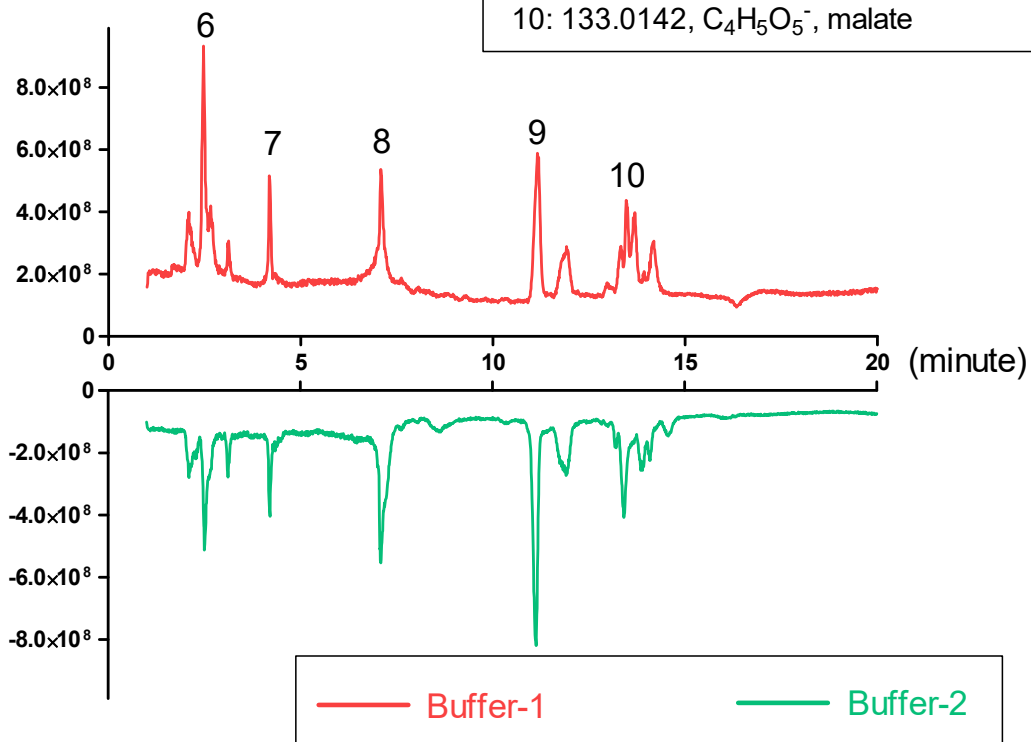


Figure S7. Annotation problems can result from two ions with similar masses that are not fully resolved. Positive mode mass spectrum at 3.13 min shows four peaks around m/z 808: m/z 806.5685 and m/z 807.5715 were correctly annotated while m/z 808.5801 and m/z 809.5870 did not match any HMDB database formulae. The peak at m/z 808.5801 actually has contributions from two overlapping peaks that are not fully resolved at a resolving power of 80,000: the $^{13}C_2^{12}C_{44}H_{81}NO_8P^+$ peak at m/z 808.5759 (red trace), and the $^{12}C_{46}H_{83}NO_8P^+$ peak at m/z 808.5851 (green trace).

(A) liver, positive mode



(B) liver, negative mode



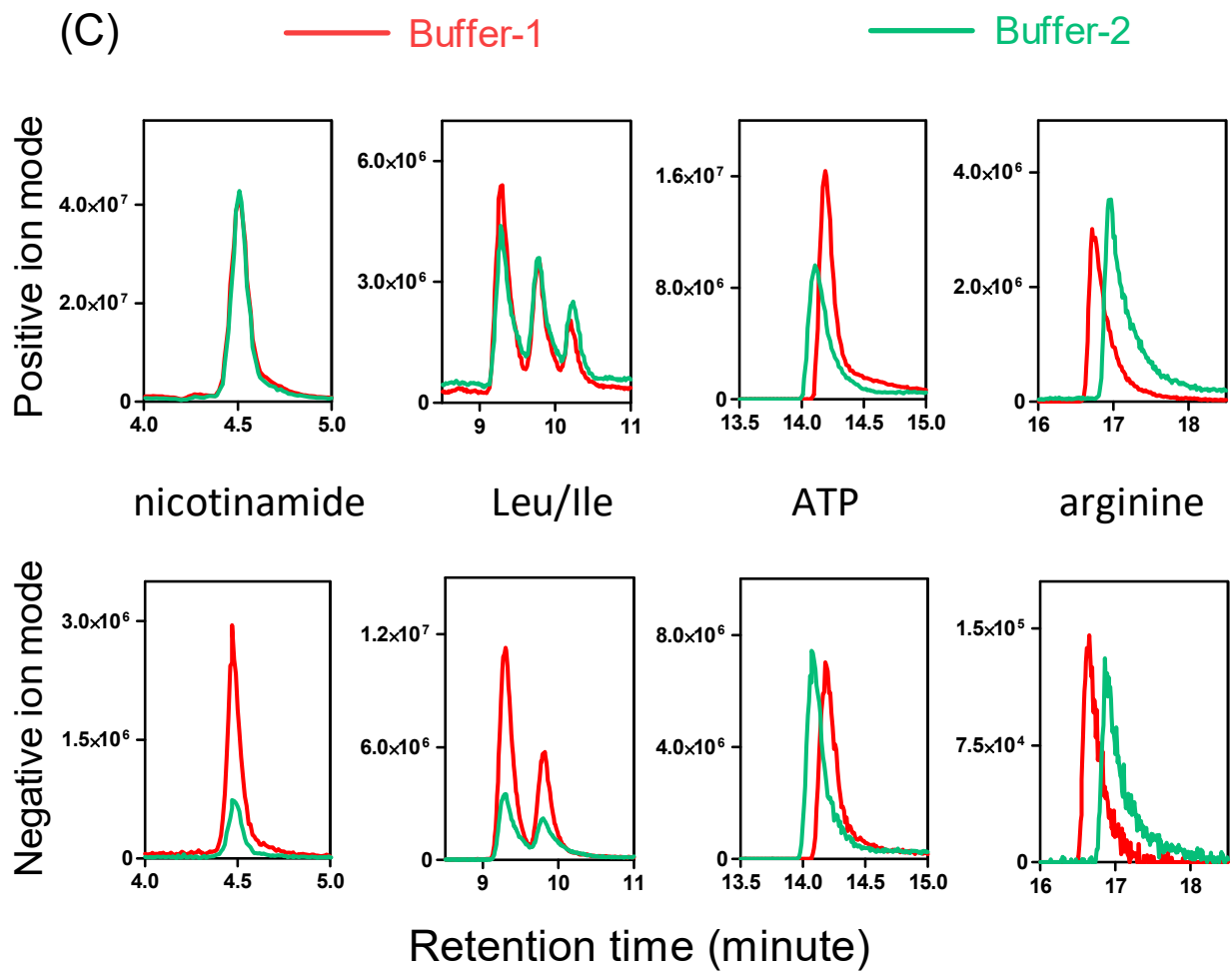


Fig S8. Total ion chromatogram (TIC) of liver extract and extracted ion chromatograms (EIC) of selected metabolites in Buffer-1 and Buffer-2. (A) Liver TIC in positive mode. At selected RTs, the most abundant species is indicated with information on m/z , formula, and annotation. (B) liver TIC in negative mode. (C) EICs of nicotinamide, leucine (9.3 min)/isoleucine (9.8 min), ATP, and arginine.

Table S2. BMW workflow for *S. cerevisiae* data analysis, showing the datasets used for different annotation steps.

Steps		Notes	Buffer-1			Buffer-2
			OD=0.02 Extract at 0 eV	OD=0.80 extract at 0 eV	OD=0.80 extract at 5/10 eV	OD=0.80 extract at 0 eV
Peak extraction		Finding all peaks in <i>S. cerevisiae</i> extract (<i>m/z</i> , RT, intensity)		X		
Annotation	1	“Background”: peaks in OD=0.80 extract with intensity < 2-fold of that in OD=0.02 extract	X	X		
	2a	“FTMS artifact”: ringing peaks around strong intensity ions		X		
	2b	“Isotope”		X		
	3a	“Adduct”: other than NH ₄ ⁺ and CH ₃ COO ⁻ adducts		X		
	3b	NH ₄ ⁺ and CH ₃ COO ⁻ adducts		X		X
	3c	“Buffer sensitive”: peaks with large intensity change when switching buffer		X		X
	3d	“Multicharge”/“Dimer”		X		
	4	“Fragment”: peaks with intensity increasing significantly at 5 or 10 eV of in-source CID, with special rules for peaks with abundant adducts		X	X	
Metabolite identification		Remaining peaks are considered “Putative metabolite” and searched against database to find formula matches		X		

Table S3. PAVE workflow for *S. cerevisiae* data analysis, showing the datasets used for different annotation steps.

Steps	Notes	Buffer-1					
		unlabeled			¹⁵ N	¹³ C	¹³ C/ ¹⁵ N
		OD=0.02 extract at 0eV	OD=0.80 extract at 0eV	OD=0.80 extract at 5/10eV	OD=0.80 extract at 0eV	OD=0.80 extract at 0eV	OD=0.80 extract at 0eV
1	Peak extraction: finding all peaks (<i>m/z</i> , RT, intensity)		X				
2	ATOMCOUNT: determining C/N counts based on mass shifts across unlabeled and labeled samples**		X		X	X	X
3	Annotating "Background": peaks in OD=0.80 extract with intensity < 2-fold of that in OD=0.02 extract*	X	X				
4	Annotating "FTMS artifact": ringing peaks around strong intensity ions*		X				
5	Annotating "Nonbio": peaks without labeling pattern **		X		X	X	X
6	Annotating "Isotope"*		X				
7	Annotating "Adduct": using additional constraint that adducts should have same C/N numbers as metabolites**		X		X	X	X
8	Annotating "Multicharge"/"Dimer": using additional constraint on C/N number**		X		X	X	X
9	Annotating "Low_C": Too low C count for mass**		X		X	X	X
10	Annotating "Fragment": peaks with intensity increasing significantly at 5 or 10 eV of in-source CID, with special rules for peaks with abundant adducts*		X	X			
11	Metabolite identification: finding putative metabolites that match database formula and C/N numbers (those without C/N matching are annotated as "Unknown")**		X		X	X	X

*: Uses the same rules as BMW

** : Uses different rules from BMW. "ATOMCOUNT" determines the C/N numbers for every peak in step 1 by examining the mass and intensity shifts between unlabeled and labeled samples. "Nonbio" are those peaks without clear labeling pattern so that C/N numbers cannot be determined. "Adduct"/"Multicharge"/"Dimer" use same constraints on $\Delta m/z$, RT and intensity ratio, with additional constraint on C/N number, e.g. adducts must have same C/N numbers as the corresponding metabolite peaks.

Table S7. Summary of liver LC-MS peak annotation*

Category and subcategory	Positive mode							Negative mode						
	>10 ³	>10 ⁴	>10 ⁵	>10 ⁶	>10 ⁷	>10 ⁸	Total	>10 ³	>10 ⁴	>10 ⁵	>10 ⁶	>10 ⁷	>10 ⁸	Total
Total	2379	8995	4109	847	152	26	16508	1027	6546	2254	423	52	6	10309
Background (1)	989	2449	757	84	10		4289	282	1354	350	47	10		2043
FTMS artifact (2a)	77	284	92	2			455	55	319	59	2			435
Isotope (2b)	205	1424	887	190	30	6	2742	96	1068	462	84	4		1714
Adduct (3a/3b)	109	759	453	114	11		1446	101	930	337	47	2		1417
Buffer sensitive (3c)	346	992	248	40	1		1627	175	603	147	41	5	1	973
Multicharge/Dimer (3d)	1	45	30	4	2		82	1	72	54	10	1		138
Fragment (4)	54	360	187	36	1		638	31	230	123	31	2	0	417
Reaction product (5)	87	475	379	81	10		1032	24	229	88	21	3	1	366
Putative metabolite without formula match	396	1634	641	81	20	3	2775	182	1069	318	40	2	0	1611
Putative metabolite with formula match	115	573	435	215	67	17	1422	80	672	316	100	23	4	1195

*: >10³ in second row is for any peak with intensity between 10³ and 10⁴, and so on