

1 **Supplementary Information**

2 **Systems-Level Annotation of a Metabolomics Data Set Reduces**
3 **25,000 Features to Fewer than 1,000 Unique Metabolites**

4

5 Nathaniel G. Mahieu and Gary J. Patti*

6

7

8

9

10

11

12

13

14

15 Department of Chemistry, Washington University, St. Louis, Missouri 63130, United
16 States

17

18

19 *Contact: gjpattij@wustl.edu, 314-935-3512

20

21 **Table S1.** Feature counts detected with various XCMS parameters.

Feature Count	method	ppm	peakwidth	prefilter	snthresh
56,368	"centWave"	6	c(2, 8)	c(4, 1000)	10
36,523	"centWave"	6	c(4, 8)	c(4, 1000)	20

22

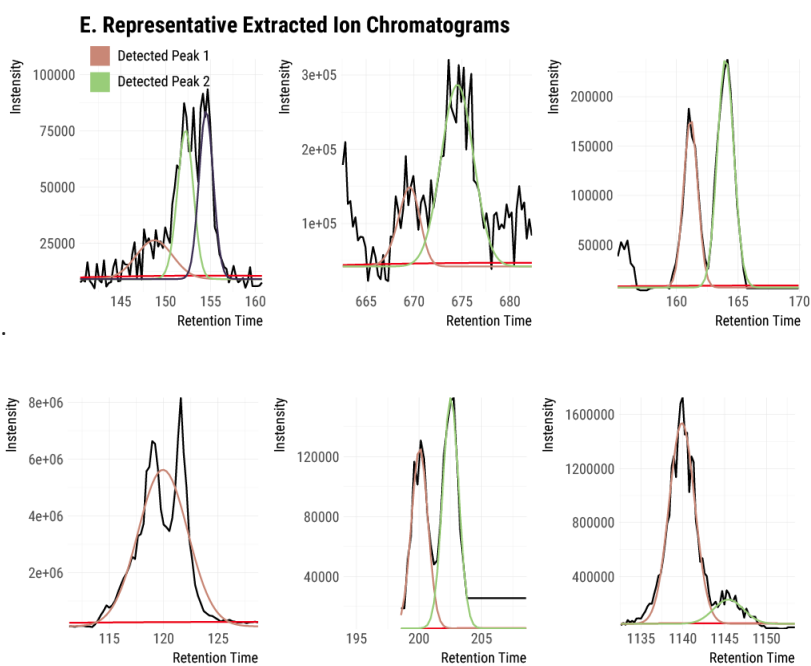
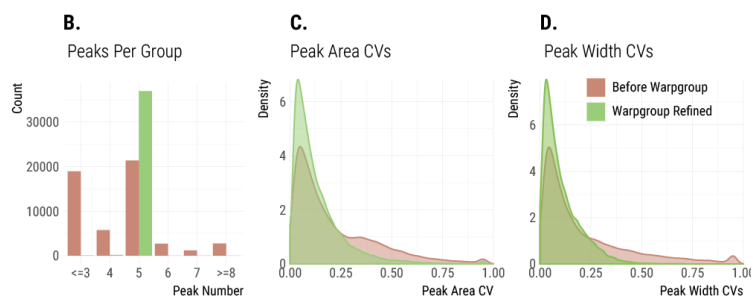
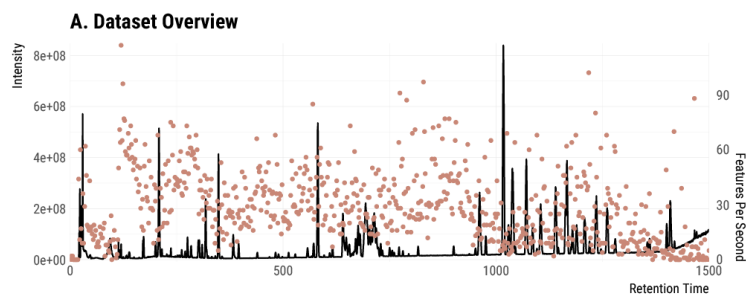
23 **Table S2.** A breakdown of the analyte number observed after each annotation step.

Stage	Groups with more than one feature		Singlets	
	All Features	Credentialed Features	All Features	Credentialed Features
Blank Subtracted	0	0	12797	2462
Isotopes	3986	1066	5071	1326
Charge Carriers	3620	1137	4384	992
Neutral Losses	3640	1174	3678	790
Multimers	3400	1117	3381	712
Commons n>200	2809	1063	2472	495
Commons n>50	2149	864	1620	353
Background	1673	659	1288	233

24

25 **Figure S1.** An overview of the consensus data set. (A) The base peak chromatogram of a
 26 representative run. The number of features detected during each second is overlaid. (B) The
 27 number of features detected in each group before (pink) and after (green) Warpgroup.
 28 Inconsistencies are resolved by Warpgroup. (C) The within group CVs of peak areas is
 29 decreased by Warpgroup. (D) The within group CVs of peak width are decreased by
 30 Warpgroup. (E) Several representative features detected by the informatic workflow. The
 31 estimated baseline is plotted in red.

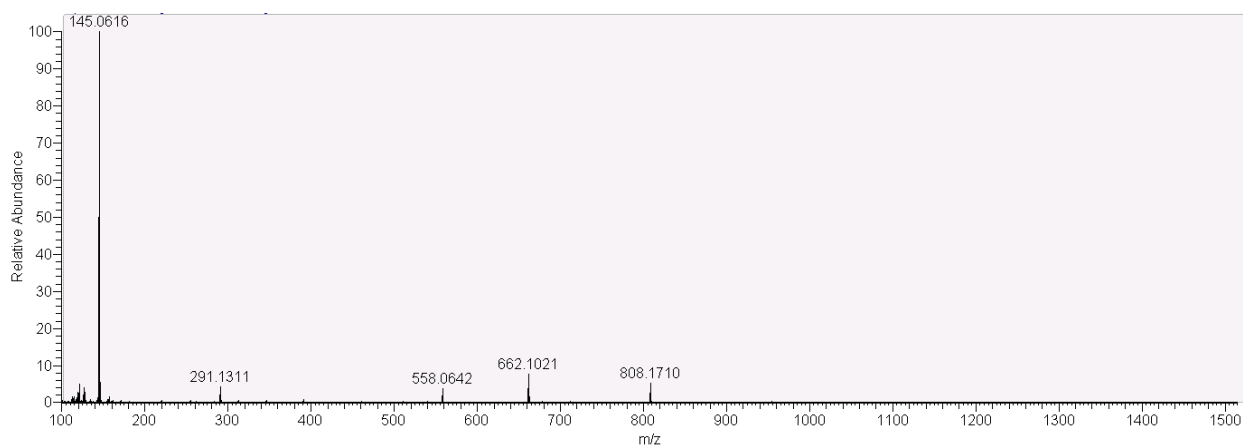
32



36
37

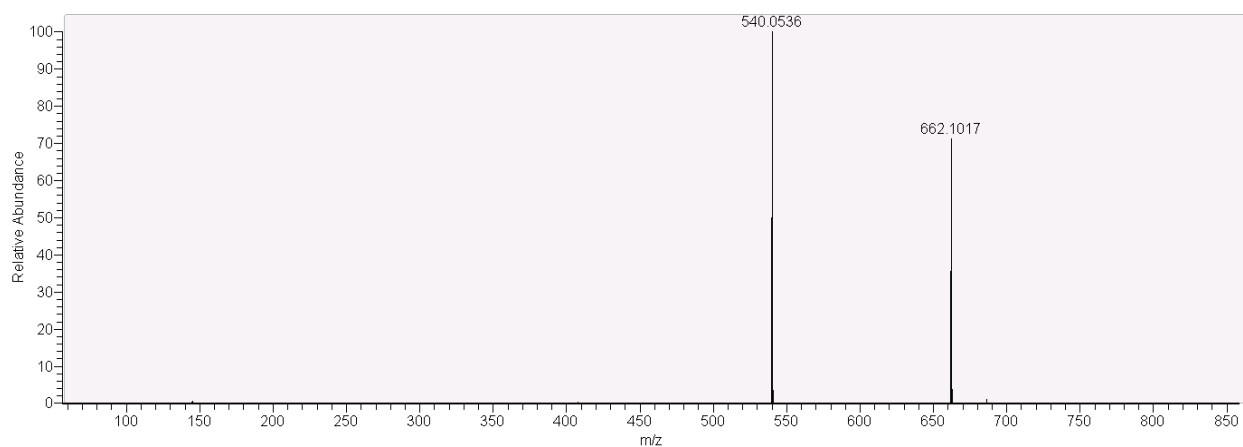
38 **Figure S2. Observation of NAD, glutamate, and their situational adduct in negative mode.**

39 A. MS1 of glutamate and NAD (negative mode electrospray ionization).



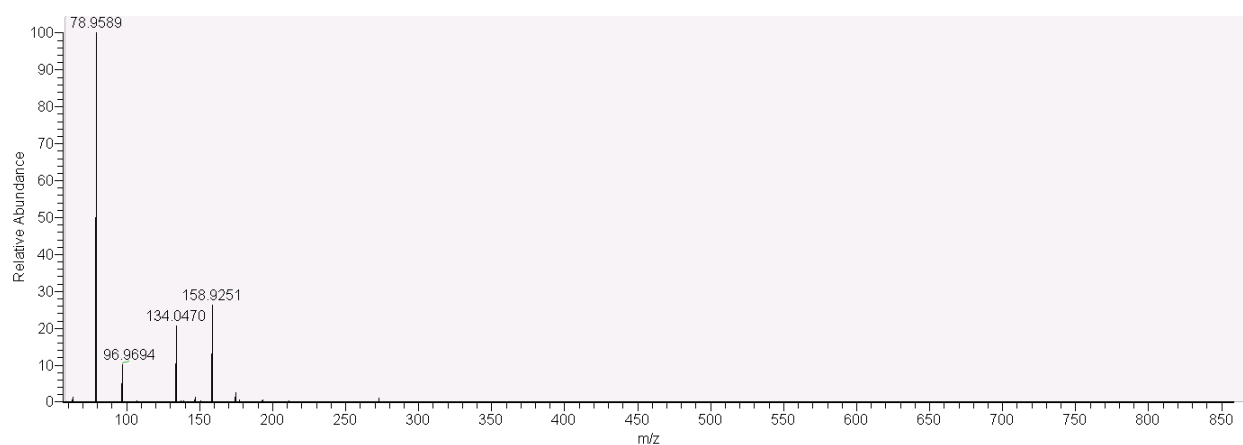
40

41 B. MS/MS of situational adduct [Glutamate + NAD - H]¹⁻, [808.1710]¹⁻. HCD = 10V.



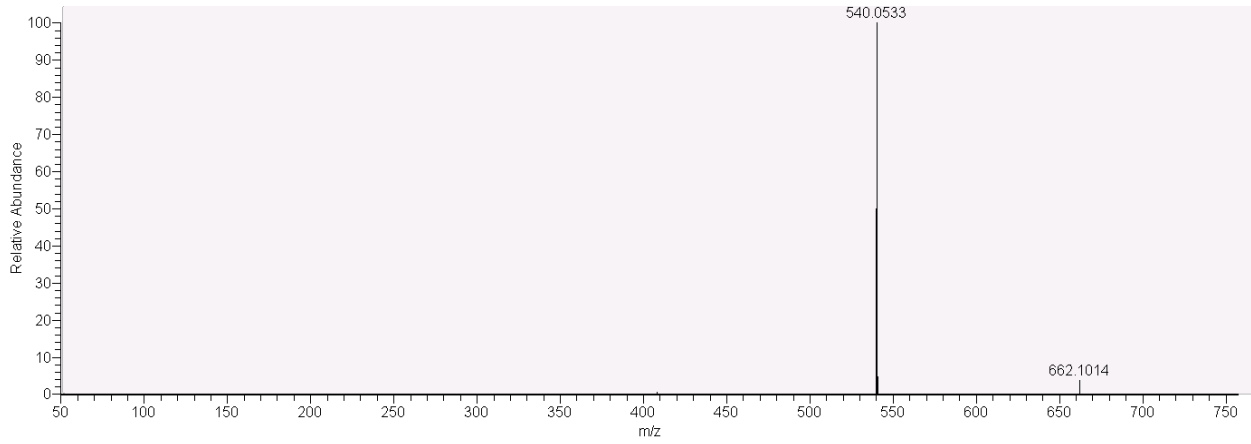
42

43 C. MS/MS of situational adduct [Glutamate + NAD - H]¹⁻, [808.1710]¹⁻. HCD = 60V.



44

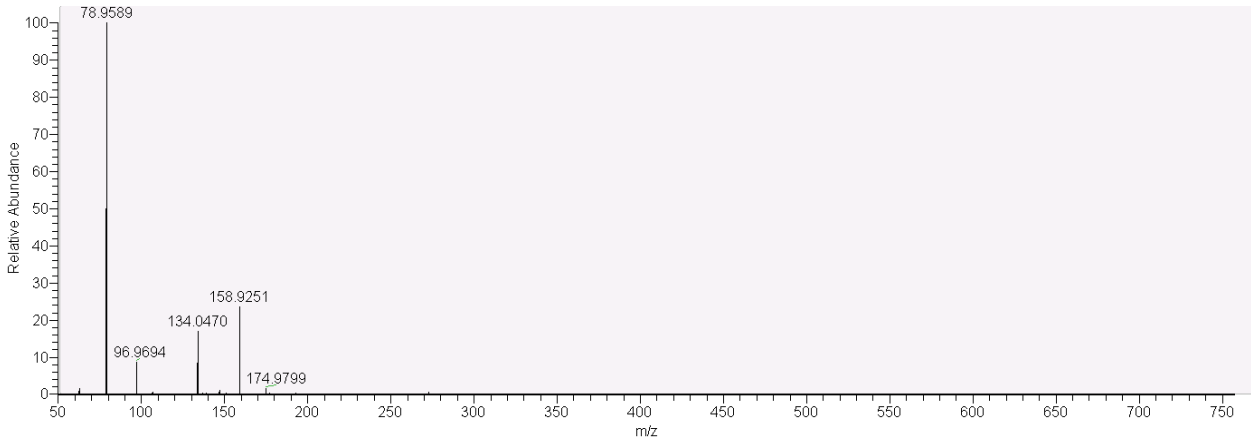
45 D. MS/MS of analyte [NAD - H]¹⁻, [662.1021]¹⁻. HCD = 10V.



46

47

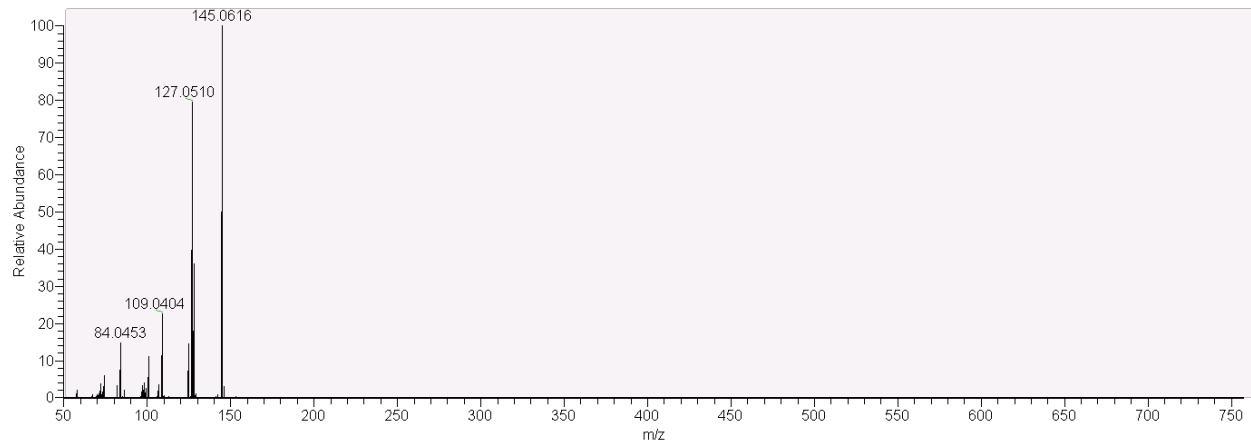
E. MS/MS of analyte [NAD - H]¹⁻, [662.1021]¹⁻. HCD = 60V.



48

49

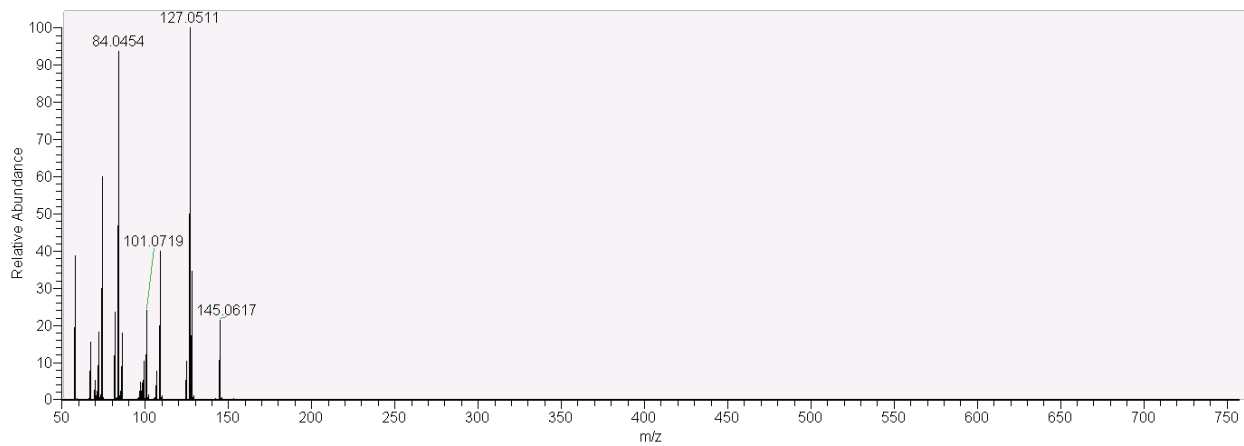
F. MS/MS of analyte [Glutamate - H]¹⁻, [145.0616]¹⁻. HCD = 10V.



50

51

G. MS/MS of analyte [Glutamate - H]¹⁻, [145.0616]¹⁻. HCD = 10V.

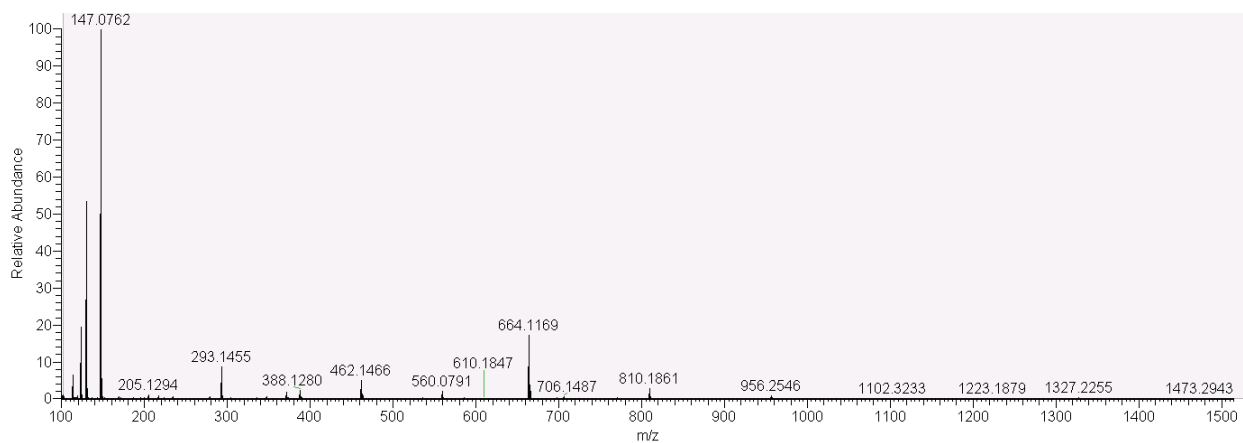


52

53

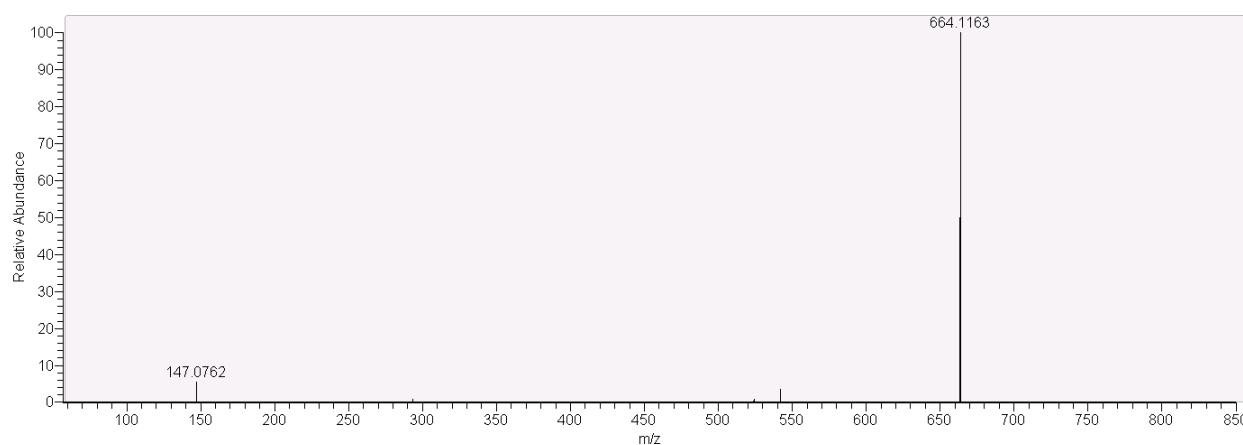
54 **Figure S3. Observation of NAD, glutamate, and their situational adduct in positive mode.**

55 **A. MS1 of glutamate and NAD (positive mode electrospray ionization).**



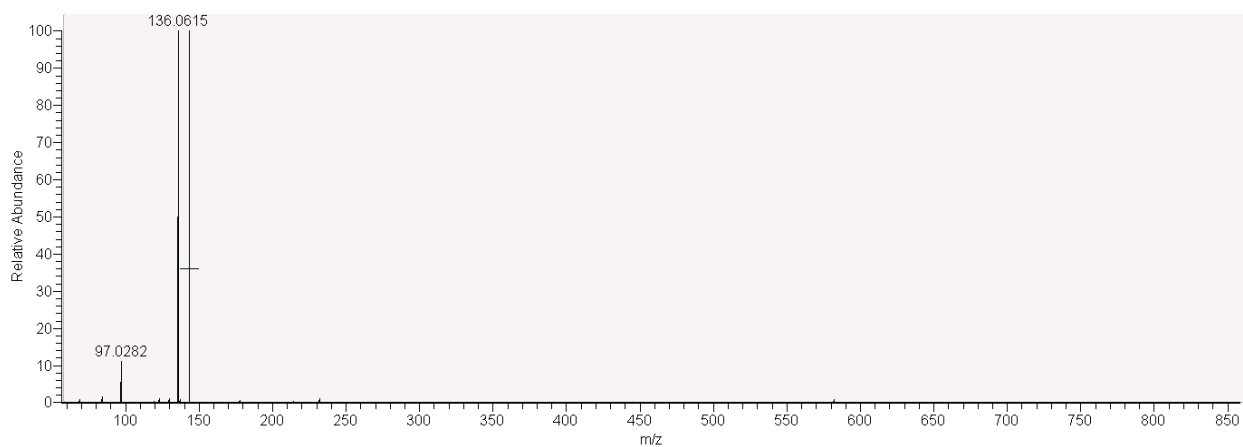
56

57 **B. MS/MS of situational adduct [Glutamate + NAD + H]¹⁺, [810.1861]¹⁺. HCD = 10V.**



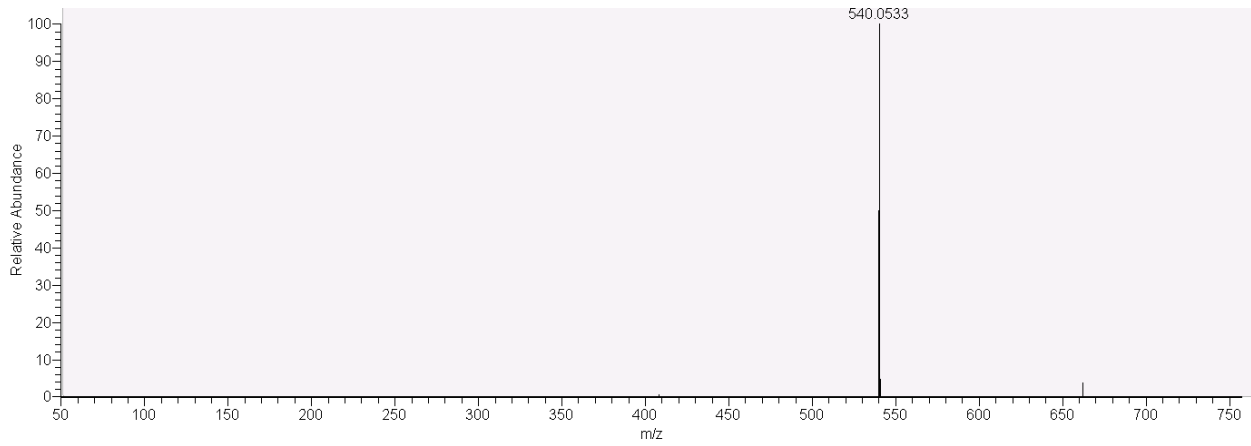
58

59 **C. MS/MS of situational adduct [Glutamate + NAD + H]¹⁺, [810.1861]¹⁺. HCD = 60V.**



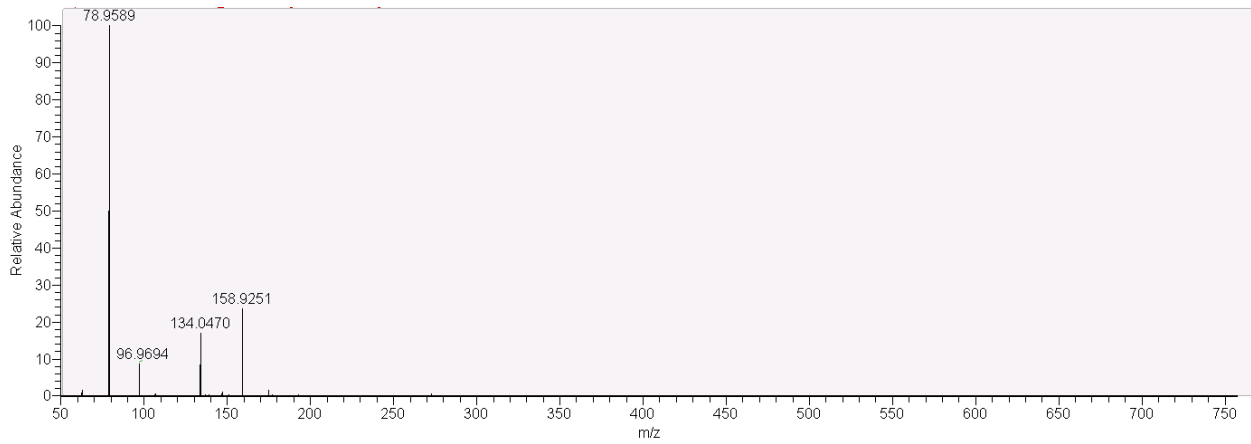
60

61 **D. MS/MS of analyte [NAD + H]¹⁺, [664.1169]¹⁺. HCD = 10V.**



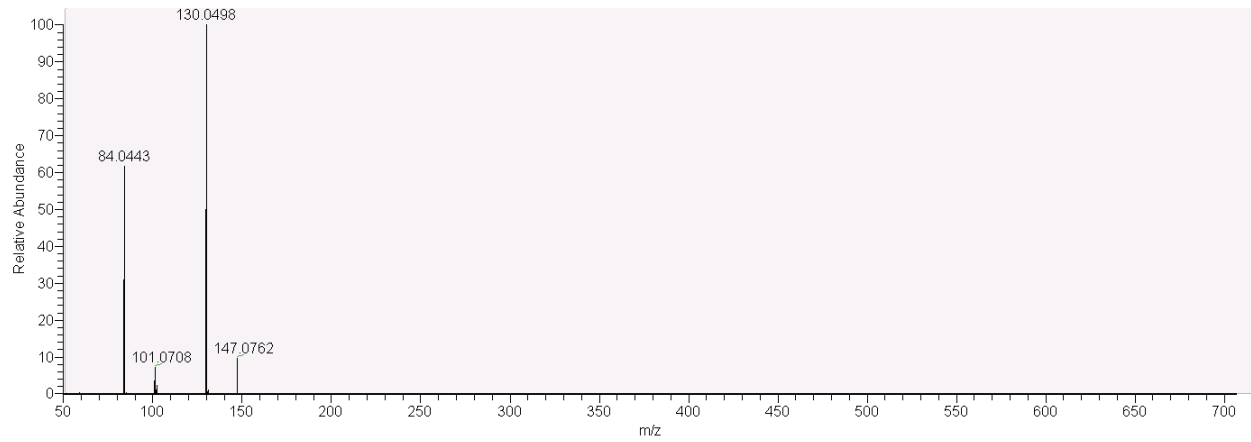
62

63 E. MS/MS of analyte $[\text{NAD} + \text{H}]^{1+}$, $[664.1169]^{1+}$. HCD = 60V.



64

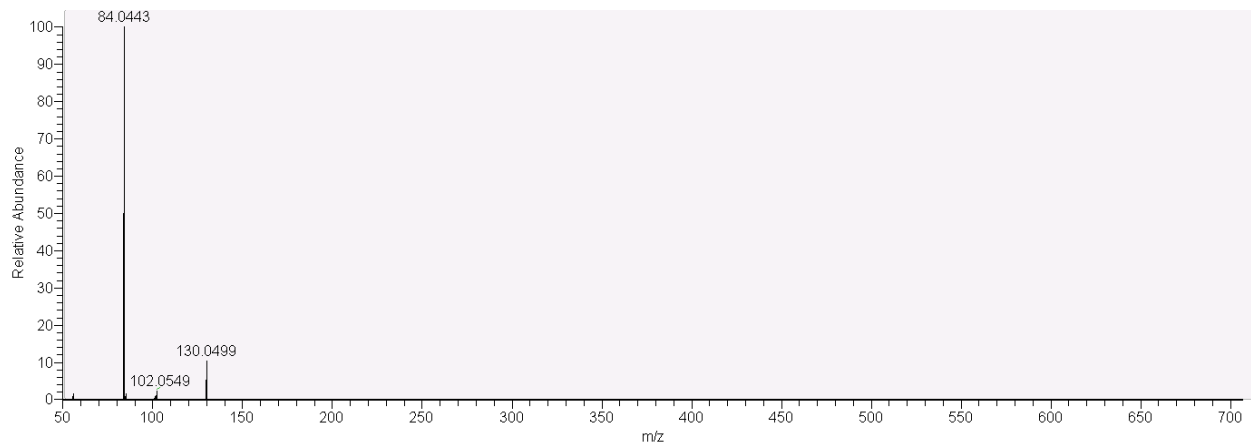
65 C. MS/MS of analyte $[\text{Glutamate} + \text{H}]^{1+}$, $[147.0762]^{1+}$. HCD = 10V.



66

67 C. MS/MS of analyte $[\text{Glutamate} + \text{H}]^{1+}$, $[147.0762]^{1+}$. HCD = 60V.

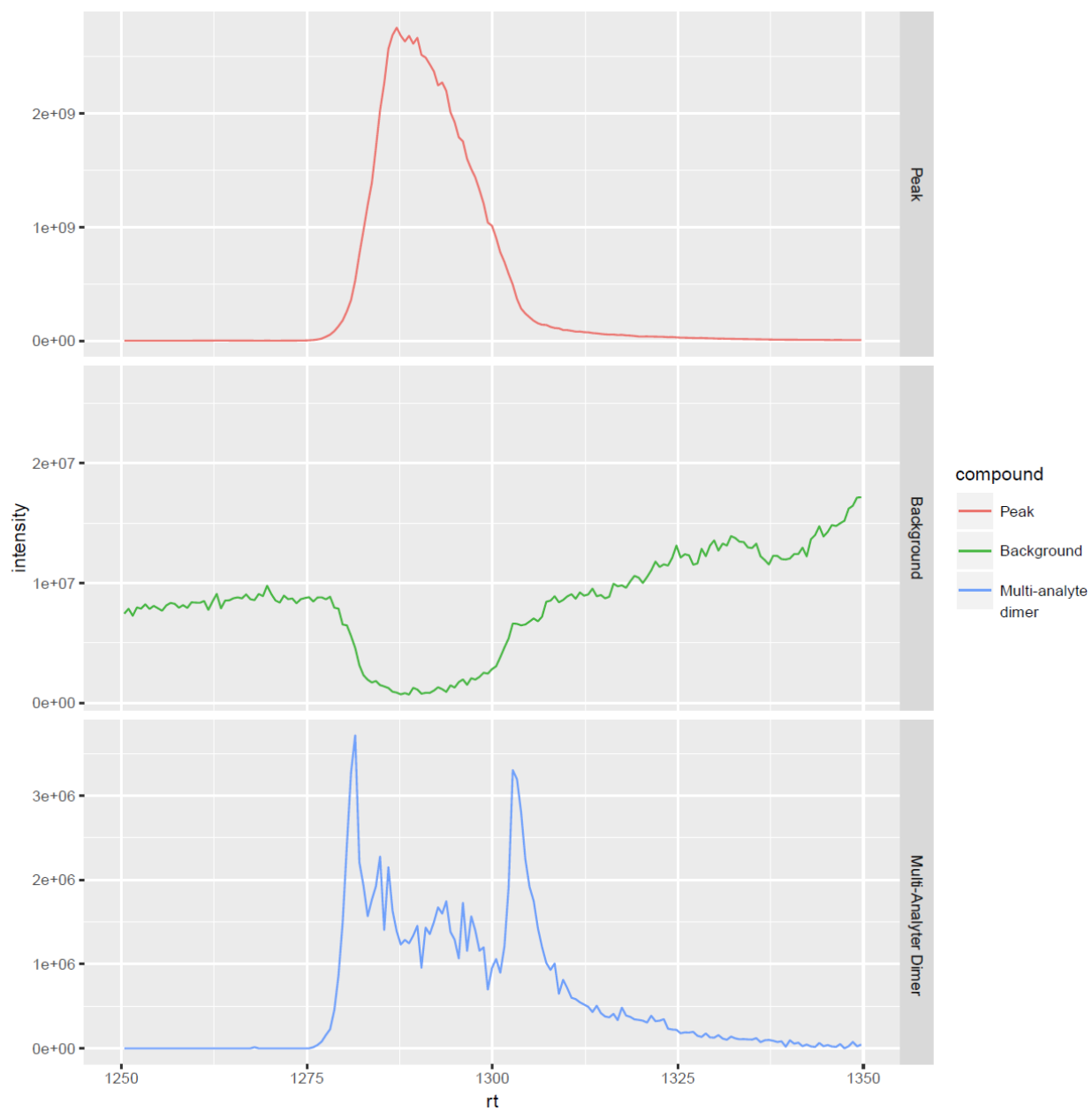
68



69

70

71 **Figure S4. Example of coeluting ions that form a multi-analyte dimer.**



72

73

74 **Figure S5.** Screenshots from the creDBle database. (A) Partial list of credentialed features
 75 showing *m/z*, retention time, polarity, grouping, and intensity. (B) A credentialed features page
 76 showing the extracted ion chromatogram, credentialed isotopes, and fragmentation data.

A.

Features

Show entries

Search:

Feature ID	Component Group	Identity	+/-	<i>m/z</i>	RT (s)	Carbons	MS/MS	Intensity
cp.g2yd14bA	cc.gA6t3nq		+	754.3843	1103.3	16	0	1.0e+07
cp.g2yd14bq	cc.gA6tmm1		+	752.3669	1057.7	0	1	2.2e+05
cp.g2yd14br	cc.gA6tmm3		+	752.4848	829.9	0	0	1.2e+05
cp.g2yd14bs	cc.gA6tmmq		+	752.3904	1091.2	0	0	2.5e+06
cp.g2yd14bt			+	753.3996	1050.7	0	0	9.8e+04
cp.g2yd14bv	cc.gA6tmmY		+	755.2841	1361.0	43	0	6.5e+05

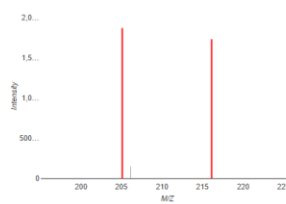
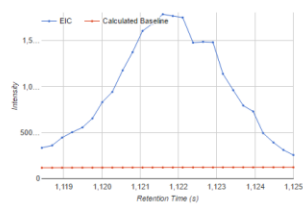
B.

cp.g2yd1mbx Details

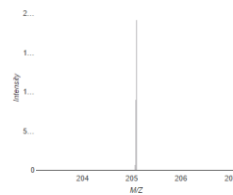
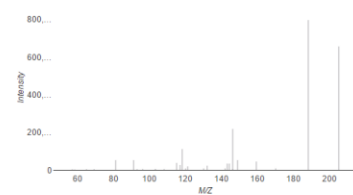
<< Previous Feature - Next Feature >>

Peak ID	<i>m/z</i>	RT	Intensity	Carbon Number	Polarity	Credentialed?	Peaks In Group	Experiment ID
cp.g2yd1mbx	710.4991	1121.6	1.7e+06	11	+	true	2	cm.h2

Spectral Data



Fragmentation Data



79 **Supplemental Methods**

80 **1 Materials**

81 U-¹³C-D-glucose was purchased from Cambridge Isotope Laboratories Inc. (Andover, MA). *E.*
82 *coli* strain K12, MG1655 was purchased from ATCC (Manassas, VA). Lennox LB broth powder
83 and 5x M9 salts were purchased from Sigma-Aldrich (St. Louis, MO). Cell culture was
84 performed with ultrapure water provided by a Milli-Q system (Millipore). LC/MS grade, Burdick &
85 Jackson brand water, acetonitrile, methanol, and isopropanol were purchased from Honeywell
86 (Morris Plains, NJ). Cortecs T3 reversed phase UPLC columns and column guards were
87 purchased from Waters Corporation (Milford, MA).

88 **Generating Credentialed Samples**

89 *E. coli* was grown in a rotary shaker at 37 °C and 300 rpm as previously described.²⁷ A 100 mL
90 volume of M9 minimal media was used with a glucose concentration of 2 g/L. Two cultures
91 were grown in parallel, one using natural-abundance glucose and a second using U-¹³C-glucose
92 as the only carbon source. Cultures were grown to OD₆₀₀ = 0.7, at which point they were
93 harvested.

94 For harvest, flasks were removed from the shaker and placed on ice. The contents of
95 each flask were pipetted into 50 mL conical tubes and centrifuged at 3200g and 4 °C for 10
96 minutes. The supernatant was decanted and remaining media was gently rinsed off the top of
97 the pellet with 0.5 mL of water. Conical tubes were then placed in liquid nitrogen and lyophilized
98 for 24 hours, or until dry. This powdered, credentialed *E. coli* standard was then extracted to
99 generate samples for untargeted metabolomic analysis.

100 Several replicate extractions were performed in parallel by using a previously described
101 method.³⁰ Briefly, five 2.5 mg samples of each ¹²C and ¹³C material were weighed out, while two
102 empty tubes were included as extraction blanks. To these, 1,000 μL of 2:2:1

103 methanol:acetonitrile:water was added, followed by three freeze-thaw cycles with sonication and
104 vortexing. After centrifugation, the supernatant was vacuum concentrated and reconstituted in
105 100 μL of 1:1 acetonitrile:water with internal standards. From these extracts, three samples
106 were aliquoted for LC/MS analysis: natural-abundance extract, a mix of 1:1 natural-abundance
107 extract and ^{13}C extract, and the blank extract.

108 ***Data Set Generation***

109 The untargeted LC/MS data set was generated in positive polarity on a Q Exactive Plus mass
110 spectrometer with a HESI II source coupled to a Dionex 3000RSLC. The data set was collected
111 with the following settings: aux gas, 5; sheath gas, 35; sweep gas, 2; capillary temperature, 300
112 $^{\circ}\text{C}$; aux gas temperature, 200 $^{\circ}\text{C}$; spray voltage, 3.5 kV; needle diameter, 34 ga; s-lens, 75 V;
113 mass range, 100–1500 Da; resolution 70,000; micro scans, 1; max injection time; 100 ms;
114 automatic gain control target: 1e6. Reversed-phase chromatography was performed with the
115 Waters Cortecs T3 (2.1mm x 50mm, 1.6 μm) column at a flow rate of 300 $\mu\text{L}/\text{min}$ and a column
116 temperature of 50 $^{\circ}\text{C}$. Solvents were: A, water + 5 mM ammonium acetate + 5 μM ammonium
117 phosphate; B, 9:1 isopropanol:methanol + 5 mM ammonium acetate + 5 μM ammonium
118 phosphate. An injection volume of 2 μL was used with a linear gradient of (minutes, %A): 0, 100;
119 28, 0; 30, 0; 30, 100; 35, 100.

120 Chromatographic features were detected. Mass traces were retained if they were longer
121 than 10 scans, excluding missing peaks. Baselines for each mass trace were calculated by
122 using the iterative restricted least squares method from the baseline R package. Model based
123 peak detection was performed by using the skew normal distribution as a model peak
124 distribution. This process resulted in a set of features detected in each replicate run. Features
125 were grouped by mass and retention time using a density based method. Retention time drift
126 and mass drift were corrected by fitting a loess curve of degree 2 to the distance from the mean
127 value of each group against the mean retention time of each group. An outline of the workflow

128 can be found on GitHub at [https://github.com/nathaniel-mahieu/metabolomic-feature-reduction-](https://github.com/nathaniel-mahieu/metabolomic-feature-reduction-2017)
129 2017.

130 Subtle variations from run to run cause many features to be integrated differently and
131 sometimes not integrated in each file. Further, closely eluting peaks often lead to incorrectly
132 grouped features. To resolve these missing values, refine the individual data sets, and get a set
133 of detected peaks consistent with all replicate runs, we applied the Warpgroup algorithm.²⁶
134 Warpgroup is available at <https://github.com/nathaniel-mahieu/warpgroup>. Warpgroup takes as
135 input the raw data and each file's detected features combining them to output a set of
136 consensus features. Parameters: sc.aligned.lim, 9; pct.pad, 0.1; min.peaks, 3. Of the detected
137 peaks we retained only features with a signal-to-noise ratio >5 and a coefficient of variation <0.5
138 after Warpgrouping. This resulted in 25,230 "high-quality" features in our representative data
139 set.

140 ***Annotation Notes***

141 Features can belong to only one group. Most group assignments are clear (e.g., a sodium
142 adduct of glutamate belongs in the glutamate group). The only ambiguous group assignments
143 are when two different metabolites dimerize (e.g., a dimer between glutamate and NAD). In this
144 case, our assignment of the dimer feature to either the glutamate or NAD group is arbitrary. We
145 note that the number of features in a group does not affect our overarching goal to count the
146 number of unique metabolites present in the data set.

147 Relationship graphs are nonlinear, span many peaks, and specify only transformations of
148 masses. M+H, M+Na, and M+K form a fully connected graph. Determination of the original
149 monoisotopic mass is a related but distinct question not answerable with current methods.