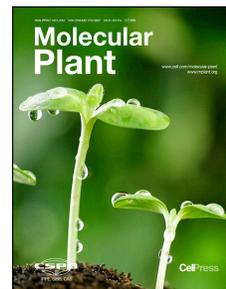


# Accepted Manuscript

Discrimination and quantification of true biological signals in LC-MS-based metabolomics analysis

Lixin Duan, István Molnár, John Hugh Snyder, Guo-an Shen, Xiaoquan Qi



PII: S1674-2052(16)30065-X  
DOI: [10.1016/j.molp.2016.05.009](https://doi.org/10.1016/j.molp.2016.05.009)  
Reference: MOLP 304

To appear in: *MOLECULAR PLANT*  
Accepted Date: 17 May 2016

Please cite this article as: **Duan L., Molnár I., Snyder J.H., Shen G.-a., and Qi X.** (2016). Discrimination and quantification of true biological signals in LC-MS-based metabolomics analysis. *Mol. Plant.* doi: 10.1016/j.molp.2016.05.009.

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

All studies published in *MOLECULAR PLANT* are embargoed until 3PM ET of the day they are published as corrected proofs on-line. Studies cannot be publicized as accepted manuscripts or uncorrected proofs.

1           **Discrimination and quantification of true biological signals in**  
2                           **LC-MS-based metabolomics analysis**

3  
4           **Running title: Filtering rules and relative calibration models for metabolomics**

5  
6  
7           Lixin Duan<sup>1</sup>, István Molnár<sup>2</sup>, John Hugh Snyder<sup>1</sup>, Guo-an Shen<sup>1</sup> and Xiaoquan Qi<sup>1\*</sup>

8  
9           <sup>1</sup>The Key Laboratory of Plant Molecular Physiology, Institute of Botany, Chinese Academy of Sciences,  
10                           20 Nanxincun, Xiangshan Road, Beijing 100093, China

11           <sup>2</sup> Natural Product Center, School of Natural Resources and the Environment, The University of Arizona,  
12                           250 East Valencia RD., Tucson, Arizona 85706, United States

13  
14           \*To whom correspondence should be addressed. E-mail: xqi@ibcas.ac.cn,  
15                           tel. +86-10-6283 6671, fax +86-10-8259 9701

20 **Dear Editor,**

21 Metabolomics is a rapidly emerging field of post-genomic research that aims to  
22 comprehensively analyze all metabolites in biological samples. Potential biomarkers that  
23 distinguish prostate cancer samples were successfully identified through metabolomics  
24 analysis (Sreekumar et al., 2009). Metabolome quantitative trait loci (mQTL) and  
25 genome-wide association studies coupled with metabolomics analysis (mGWAS) also  
26 became efficient tools to decipher the genetic basis of complex metabolic traits in large  
27 populations (Gong et al., 2013; Chen et al., 2014).

28 Liquid chromatography-mass spectrometry (LC-MS) techniques widely used for  
29 metabolomics analysis allow for the highly sensitive, high throughput detection of  
30 thousands of metabolites (Chen et al., 2013). However, LC-MS unavoidably yields a large  
31 number of false positive signals mixed with true biological signals (Yu et al., 2013; Broadhurst  
32 and Kell, 2007). Without the means to confidently discriminate and evaluate the detected  
33 signals, biomarker discovery turns out to be misleading, or downright impossible.  
34 Metabolomics is routinely used to compare the relative concentrations of metabolites  
35 amongst different samples. For this, the monitored signals must fall within the quantitative  
36 dynamic range and show a good quantitative correlation with the amounts of the  
37 compounds of interest. Thus, it is necessary to systematically evaluate the quantitative  
38 performance of each peak. Isotope-labeling of internal standards or the whole metabolome  
39 (Giavalisco, 2009) are powerful tools to improve compound annotation and relative  
40 quantification. Standard mixtures (Phinney et al., 2013) can also be used to quantitatively  
41 analyze a selected set of known metabolites. Artificial biological gradients (Redestig et al.,  
42 2011) may allow the exploration of matrix effects and quantification performance. Here, we  
43 report a novel strategy for a comprehensive LC-MS-based metabolomics analysis that  
44 enables the unambiguous and facile discrimination of biological and non-biological signals,  
45 and improves the quantification accuracy of metabolites without labeling or other  
46 specialized techniques.

47 For a LC-MS-based metabolomics experiment, we prepare a blank sample, a Quality  
48 Control mixture (QC\_mix) that combines all samples in equal proportions, and a dilution  
49 series of the QC\_mix (Figure 1A). First, three to six replicates of the blank sample are

50 analyzed to balance the instrument. This also precludes contaminating the system with the  
51 biological samples. Next, the QC\_mix is analyzed in six or more technical replicates to  
52 identify peaks that can be reproducibly detected. Third, the dilution series of the QC\_mix is  
53 analyzed, proceeding from the most dilute (16 times dilution, DS\_1/16x) to the most  
54 concentrated (two times concentration, DS\_2X) (Figure 1A). Finally, the individual samples  
55 are analyzed separately in a random order. This data acquisition pipeline enables the  
56 discernment of true biological signals, and the building of calibration curves for the  
57 quantification of all metabolites, including unknown peaks.

58 Figure 1B illustrates the principles of distinguishing signals derived from true  
59 metabolites from those that have a non-biological origin. Peak 3 is absent from the blank  
60 and its signal intensity displays a good quantitative response in the dilution range. Such  
61 peaks are considered to originate from the biological source and have good quantitative  
62 performance. Peak 2 is detected in the blank sample and its signal intensity is independent  
63 of the dilution. Such peaks may derive from the chromatography solvent, or from column  
64 contaminants. Peak 1 is also present in the blank sample, but its signal intensity is  
65 dependent on the dilution. Such peaks may represent contaminations introduced during  
66 biological extract preparation, or impurities from labware (Supplemental Table S1).

67 We developed a hierarchical five-step filtering approach which applies these principles  
68 to comprehensive metabolomics experiments that often present thousands of peaks  
69 (Supplemental Results). To validate this LC-MS-based metabolomics strategy, we prepared  
70 and analyzed two groups of artificial samples, each including 20 standard compounds. In  
71 total, 1,342 peaks were enumerated from these artificial samples after standard peak  
72 extraction and alignment (Figure 1C). Step 1, the reproducibility check, eliminated 1,053  
73 peaks as these were detectable less than five times in the six replicates of the QC\_mix. Step  
74 2, the variation check, filtered out a further 40 peaks since their relative standard deviation  
75 (RSD) was >20% in the six QC\_mix replicates. Another 104 peaks failed to satisfy Step 3, the  
76 blank check, since they showed a peak area ratio of blank to sample (Ratio<sub>B/S</sub>) of > 1%. Step  
77 4, the response check, eliminated an additional 28 peaks because these had an  
78 unsatisfactory quantitative correlation ( $r < 0.9$ ) in the QC\_mix dilution series. Finally,  
79 re-extraction and manual inspection of the peaks with an  $r$  between 0.9 and 0.99 (Step 5)

80 disqualified a further 17 signals, resulting in a final set of 102 peaks. Remarkably, all 20  
81 standard compounds in the artificial samples passed this strict, hierarchical filtering process,  
82 while 92.4% of the peaks were eliminated as false positives or peaks with insufficient  
83 quantitative performance. Of the final set of 102 filtered peaks, 53 were identified to derive  
84 from the 20 standard compounds and their adductor fragment ions. A further 36 peaks were  
85 de-replicated to 21 unknown metabolites. However, these peaks were also deduced to  
86 originate from the 20 standards, because their concentration ratios were very similar to  
87 those of the standards in the artificial samples (Supplemental Figure S5, Supplemental  
88 information and Additional supplementary Data 1). These unknown compounds may be  
89 minor impurities present in the standards. Only 13 peaks assigned to 11 compounds had an  
90 unknown origin.

91 In a biological sample, the absolute concentration of most compounds represented by a  
92 peak will remain necessarily unknown due to the absence of standards. To overcome this  
93 limitation, our strategy introduces a relative concentration index (RCI) as an arbitrary  
94 concentration. To calculate the RCI, we build a relative quantitative model for the calibration  
95 of each analyte in the QC\_mix dilution series, using a custom Python script (Supplemental  
96 Methods). We assume that the RCI of any compound is 3,200, 1,600, 800, 400, 200, and 100  
97 arbitrary units in a 2x, 1x, 1/2x, 1/4x, 1/8x, and 1/16x sample dilution, respectively. Using  
98 these assumptions, the coumarin peak, for example, yielded a relative calibration model of  $y$   
99  $= 435.07x + 24,301$  for the dilution series (Supplemental Figure S1A). Since the peak area of  
100 coumarin in one of the artificial samples was 672,060 units, its RCI in that sample was  
101 calculated to be 1,489 arbitrary units.

102 Among the 102 filtered peaks in our validation experiment, 72.0% fitted a linear model,  
103 21.8% fitted a binomial model, and 3.4% fitted a logarithmic model. A combination of a  
104 linear and a binomial model was manually built for one particular peak (Supplemental Figure  
105 S1D). To evaluate the quantitative accuracy of our strategy, we compared the RSDs of the  
106 ratios of the concentrations to peak areas on one hand and to RCIs on the other  
107 (Supplemental Table S2). Indeed, for 75% of the standards (15 out of 20), the quantification  
108 was more reliable when using the RCI. In addition, because all peaks are calibrated using the  
109 same scale range, we can use RCI instead of the peak area to compare changes amongst

110 various samples in a more precise manner (Supplemental Results). Moreover, by following  
111 our metabolomics analysis strategy, metabolite peaks that pass the five-step filtering process  
112 and quantified by RCI can be used as a targeted metabolomics dataset. A metabolite report  
113 is also generated displaying the final metabolite list, with each filtered peak annotated with  
114 the unique  $m/z$ , the retention time and other evaluation parameters (Supplemental Figure  
115 S2B).

116 We further applied our strategy for the analysis of the metabolomes of the seeds of  
117 two typical rice cultivars, 9311 (*Oryza sativa* L. ssp. *indica*) and Nipponbare (*O. sativa* L. ssp.  
118 *japonica*). A total of 2,162 peaks were enumerated from all samples, but 71.3% of these  
119 initial peaks were eliminated using the five-step filtering approach (Supplemental Figure S3).  
120 Principal component analysis (PCA) of the variations between the two rice samples yielded a  
121 model with both a better explanatory performance and a higher predictive power for our  
122 strategy, compared to that for the traditional, non-targeted metabolomics analysis method  
123 (Figure 1D, F and G, and Supplemental Figure S4).

124 Our new strategy also significantly reduced the number of false positive differential  
125 peaks, filtering out 444 of the 565 peaks that a traditional metabolomics analysis may have  
126 considered as potential biomarkers for the two rice cultivars (Figure 1E, Supplemental Table  
127 S3). For example, the area of the peak 636.2171\_0.6162 ( $m/z$ \_Rt) is 6.3 times larger in the  
128 9311 cultivar sample than in the Nipponbare sample. This peak satisfied the reproducibility  
129 and the variation criteria, and was absent in the blank samples. However, it showed an  
130 unacceptable quantitative performance in the response check (Step 4), with a correlation  
131 coefficient of -0.2415 between the peak areas and the RCI. Thus, although peak  
132 636.2171\_0.6162 may be flagged as a distinguishing biomarker by a traditional  
133 metabolomics pipeline, it is revealed by our analysis to represent a compound of  
134 non-biological origin. Even more remarkably, the metabolite report containing our final list  
135 of differential peaks allowed the putative identification of 30 metabolites by database  
136 comparisons, including 12 lipids, nine flavonoids, two amino acids, two phenolics, two  
137 nucleosides, vitamin B6, hydroxylamine and a diterpenoid (Supplemental Table S4). Many of  
138 these are validated biomarkers critically important for the yield and the seed quality of the  
139 two rice cultivars (Supplemental Results). Nevertheless, possible interactions between the

140 biological matrix and some compounds in a complex metabolome may complicate the  
141 identification of filtered biomarkers.

142 In conclusion, our new strategy greatly reduces the number of false positive peaks,  
143 enhances quantitative accuracy, and allows for a more meaningful analysis of comparative  
144 metabolomics datasets.

145

#### 146 **SUPPLEMENTAL INFORMATION**

147 Supplemental information is available at Molecular Pant Online.

#### 148 **AUTHOR CONTRIBUTIONS**

149 L.X.D., X.Q. and J.H.S designed the research. L.X.D. performed the experiments. G.A.S. wrote  
150 the script to build the optimized regression models. L.X.D., X.Q. and I.M. analyzed the data  
151 and wrote the article.

#### 152 **FUNDING**

153 This work was supported by the Chinese National Key Programs for Research and  
154 Development (JFYS2016ZY03002153), the Key Project of Chinese National Programs for  
155 Fundamental Research and Development (2013CB127000), the Strategic Priority Research  
156 Program of the Chinese Academy of Sciences (XDA08020104), the National High-tech R&D  
157 Program of China (2012AA10A304-3) and the National Natural Science Foundation of China  
158 (31200227).

#### 159 **ACKNOWLEDGMENTS**

160 No conflict of interest declared.

161

## REFERENCES

- 162  
163  
164 Broadhurst, D.I., and Kell, D.B. (2007). Statistical strategies for avoiding false discoveries in metabolomics and  
165 related experiments. *Metabolomics* 2:171-196.
- 166 Chen, W., Gao, Y., Xie, W., Gong, L., Lu, K., Wang, W., Li, Y., Liu, X., Zhang, H., Dong, H., Zhang, W., Zhang, L., Yu,  
167 S., Wang, G., Lian, X., and Luo, J. (2014). Genome-wide association analyses provide genetic and  
168 biochemical insights into natural variation in rice metabolism. *Nat Genet* 46:714-721.
- 169 Chen, W., Gong, L., Guo, Z., Wang, W., Zhang, H., Liu, X., Yu, S., Xiong, L., and Luo, J. (2013). A novel integrated  
170 method for large-scale detection, identification, and quantification of widely targeted metabolites:  
171 Application in the study of rice metabolomics. *Mol Plant* 6: 1769-1780.
- 172 Giavalisco, P., Kohl, K., Hummel, J., Seiwert, B., and Willmitzer, L. (2009). <sup>13</sup>C isotope-labeled metabolomes  
173 allowing for improved compound annotation and relative quantification in liquid  
174 chromatography-mass spectrometry-based metabolomic research. *Anal Chem* 81: 6546-51.
- 175 Gong, L., Chen, W., Gao, Y., Liu, X., Zhang, H., Xu, C., Yu, S., Zhang, Q., and Luo, J. (2013). Genetic analysis of the  
176 metabolome exemplified using a rice population. *Proc Natl Acad Sci U S A* 110:20320-20325.
- 177 Phinney, K.W., Ballihaut, G., Bedner, M., Benford, B.S., Camara, J.E., Christopher, S.J., Davis, W.C., Dodder, N.G.,  
178 Eppe, G., Lang, B.E., Long, S.E., Lowenthal, M.S., McGaw, E.A., Murphy, K.E., Nelson, B.C., Prendergast,  
179 J.L., Reiner, J.L., Rimmer, C.A., Sander, L.C., Schantz, M.M., Sharpless, K.E., Sniegoski, L.T., Tai, S.S.,  
180 Thomas, J.B., Vetter, T.W., Welch, M.J., Wise, S.A., Wood, L.J., Guthrie, W.F., Hagwood, C.R., Leigh, S.D.,  
181 Yen, J.H., Zhang, N.F., Chaudhary, W.M., Chen, H., Fazili, Z., LaVoie, D.J., McCoy, L.F., Momin, S.S.,  
182 Paladugula, N., Pendergrast, E.C., Pfeiffer, C.M., Powers, C.D., Rabinowitz, D., Rybak, M.E., Schleicher,  
183 R.L., Toombs, B.M., Xu, M., Zhang, M., and Castle, A.L.. (2013). Development of a Standard Reference  
184 Material for metabolomics research. *Anal Chem* 85:11732-11738.
- 185 Redestig, H., Kobayashi, M., Saito, K., and Kusano, M. (2011). Exploring matrix effects and quantification  
186 performance in metabolomics experiments using artificial biological gradients. *Anal Chem*  
187 83:5645-5651.
- 188 Sreekumar, A., Poisson, L.M., Rajendiran, T.M., Khan, A.P., Cao, Q., Yu, J., Laxman, B., Mehra, R., Lonigro, R.J., Li,  
189 Y., Nyati, M.K., Ahsan, A., Kalyana, S.S., Han, B., Cao, X., Byun, J., Omenn, G.S., Ghosh, D., Pennathur, S.,  
190 Alexander, D.C., Berger, A., Shuster, J.R., Wei, J.T., Varambally, S., Beecher, C., and Chinnaiyan, A.M.  
191 (2009). Metabolomic profiles delineate potential role for sarcosine in prostate cancer progression.  
192 *Nature* 457:910-914.
- 193 Yu, T., Park, Y., Li, S., and Jones, D.P. (2013). Hybrid Feature Detection and Information Accumulation Using  
194 High-Resolution LC-MS Metabolomics Data. *J Proteome Res* 12:1419-1427.
- 195

196 **Figure 1. The strategy of LC-MS-based metabolomics analysis.** (A) The design of the  
197 metabolomics experiment. DS\_Nx, a sample in the dilution series with an Nx dilution factor;  
198 RCI, relative concentration index; SIM, selected ion monitoring; MRM, multiple reaction  
199 monitoring. (B) Overlay of LC-MS total ion chromatograms for a dilution series of a sample  
200 and the blank samples. (C) Benchmarking the five steps of peak filtering using artificial  
201 samples. (D) Principle components analysis (PCA) models for rice samples. Left, the  
202 traditional method; right, the new strategy.  $R^2X$  (cum), cumulative explained variation;  $Q^2$   
203 (cum), cumulative cross validated  $R^2$ ; Comp, principal component identified in the model. (E)  
204 Venn diagram of the number of differential peaks identified in the traditional (circle on the  
205 left) versus the new metabolomics approach (circle on the right). (F) and (G) The loading  
206 S-plots for PCA with the traditional metabolomics method (F) and the new strategy using RCI  
207 (G).  
208

