# Untargeted metabolomics suffers from incomplete data analysis

**Richard Baran**

Baran Bioscience, LLC; 2150 Allston Way Suite 400; Berkeley, California, USA.

E-mail: richard@baranbioscience.com

*Brief Communication*

**Abstract**

*Introduction:* Untargeted metabolomics is a powerful tool for biological discoveries. Significant advances in computational approaches to analyzing the complex raw data have been made, yet it is not clear how exhaustive and reliable are the data analysis results.

*Objectives:* Assessment of the quality of data analysis results in untargeted metabolomics.

*Methods:* Five published untargeted metabolomics studies acquired using instruments from different manufacturers were reanalyzed.

*Results:* Omissions of at least 50 relevant compounds from original results as well as examples of representative mistakes are reported for each study.

*Conclusion:* Incomplete data analysis shows unexplored potential of current and legacy data.

1

**Introduction**

Mass spectrometry-based metabolomics is a powerful tool for the discovery of novel compounds, metabolic capabilities, and biomarkers (Patti el at. 2012; Sévin et al. 2015). Successful discoveries are dependent on the ability to reliably detect relevant signals in raw data and to correctly interpret the underlying spectral features of compounds (Kind & Fiehn 2007; Dunn et al. 2013; Scheubert et al. 2013; Baran & Northen 2013; Kind et al. 2017). The challenging complexity of the data analysis process is well recognized and computational tools facilitating the data analysis process are available (Weber et al. 2017). However, it is not clear how exhaustive and reliable are the current data analysis results. The quality of the results is important not only in the context of exploratory research but even more more so in the context of a strengthening trend towards large scale integration of multi-omic datasets (Perez-Riverol et al. 2017). Public repositories of metabolomics data, such as the UCSD Metabolomics Workbench (Sud et al. 2016) or the MetaboLights (Haug et al. 2013) database, provide an opportunity to reanalyze published raw data to assess the coverage of relevant signals as well as the quality of mass spectra interpretation.

Five untargeted metabolomics datasets from public repositories acquired using instruments from different manufacturers were selected for reanalysis (Table 1, Supplementary Fig. 1-5). The selection was arbitrary with a focus leaning towards high complexity of the raw data (large numbers of detected compounds).

**Materials and Methods**

Raw datafiles along with accompanying data analysis results were downloaded from the respective data repositories (Table 1). Raw data files in original instrument manufacturers' proprietary data formats were converted to mzXML (Pedrioli et al. 2004) data format using ProteoWizard's msconvert tool (Chambers et al. 2012). Differences among datasets within a specific study (for ions not reported in original study results) were detected using direct comparisons between datasets binned along the m/z dimension as described previously (Baran et al. 2006). The

2

43    mass spectra and extracted ion chromatograms corresponding to candidate differences were then

44    inspected visually to assign related ions (e.g. $[M+H]^+$, adducts, multimers, in-source fragments,

45    isotopic peaks). To limit the extent of tedious manual curation, the aim of the reanalysis was to find

46    50 relevant omissions in each study.

47        To be considered an omission, none of the ions corresponding to the omitted compound could be

48    reported in the original results (even if the only reported ion corresponds to an isotopic peak of an

49    in-source fragment ion of a specific compound). Only raw data acquired in positive mode polarity

50    were used for re-analysis for each study. However, negative mode raw data and results were

51    examined as well. If none of the ions of a specific compound were reported in positive mode

52    results, but at least one ion related to the compound was reported in negative mode results, the

53    compound was not considered and not reported as an omission.

54        Multiple ions for omitted compounds along with their peak areas are listed in Supplementary

55    Data 1. These lists of ions are not exhaustive. Low intensity isotopic peaks or ions that could be

56    potentially related (but not showing clear similarities in chromatographic profiles, relative peak

57    areas across samples, or differences in $m/z$ to other ions of typical chemical relationships) may have

58    been left out of these lists. However, records for even these possibly related ions were sought in

59    original results accompanying the study to make the best effort to report truly omitted compounds

60    in reanalysis results.

61        Peak areas were calculated using the trapezoidal integration method without any prior smoothing

62    of extracted ion chromatograms or baseline subtraction. Integration bounds were set manually. The

63    ion with the largest peak area from a group of related ions was selected as a "representative" ion for

64    a given compound and used for extracted ion chromatograms (Fig. 1, Supplementary Fig. 6-10).

65    Few representative mistakes found during the reanalysis process were mostly related to ion type

66    (mis)interpretation in the original results and are shown in Supplementary Figures 12-16. A rough

67    comparison of relevance of omitted compounds to the original results was based on peak areas of

3

68 "representative" ions and a measure of a statistical significance of a difference among the groups of

69 replicate samples in a study, if applicable (Supplementary Fig. 11). Peak areas calculated by the

70 trapezoidal method were normalized to peak areas in the original results (Supplementary Fig. 17-

71 21) for this comparison.

72 **Results and Discussion**

73     The raw data were reanalyzed as described in the Materials and Methods section to look for

74 omissions of relevant compounds as well as examples of common mistakes in the original data

75 analysis results accompanying the study data. To limit the extent of tedious manual curation of the

76 data, a goal of finding 50 relevant omissions in each study was set. For a compound to be

77 considered omitted, none of its ions (e.g. $[M+H]^+$, adducts, multimers, in-source fragments, isotopic

78 peaks) could be reported in the original results. Figure 1a-d shows a few examples of omissions

79 from one of the reanalyzed studies, and Supplementary Figures 6-10 show examples of at least 50

80 omissions from each study. These omissions are relevant in the context of reported results, since

81 these compounds show either intense signals or differ significantly among the study groups

82 (Supplementary Fig. 11). In addition to omissions, mistakes in ion type interpretation were also

83 found during the reanalysis. The most commonly observed mistake was the reporting of in-source

84 fragment ions, isotopic peaks, or other ion types instead of the protonated molecule $[M+H]^+$ ion

85 (Fig. 1e, Supplementary Fig 12-16).

86     This reanalysis of published metabolomics studies was far from exhaustive. The newly reported

87 lists of ions for omitted compounds (Suppementary Data 1) are incomplete, may contain mistakes

88 as well, and additional unreported compounds are very likely present in the raw data. The selected

89 metabolomics studies have impressive quality of the raw data as well as original data analysis

90 results which must have required significant effort and insight. And yet the results of this simple

91 reanalysis point to an additional unexplored potential of current as well as legacy metabolomics

92 data. Hopefully, these results will strengthen the appreciation for the complexities of the data

93 analysis process and further motivate improvements in computational tools and knowledgebases for

94 metabolomics data analysis.

95 **Conflict of interest**

96 The author's company Baran Bioscience, LLC provides data analysis services for metabolomics

97 and small molecule mass spectrometry.

98 **References**

99 Baran, R., & Northen, T. R. (2013). Robust automated mass spectra interpretation and chemical

100 formula calculation using mixed integer linear programming. *Analytical chemistry, 85*(20), 9777-

101 9784.

102 Baran, R., Kochi, H., Saito, N., Suematsu, M., Soga, T., Nishioka, T., et al. (2006). MathDAMP: a

103 package for differential analysis of metabolite profiles. *BMC Bioinformatics, 7,* 530.

104 Chambers, M. C., Maclean, B., Burke, R., Amodei, D., Ruderman, D. L., Neumann, S., et al.

105 (2012). A cross-platform toolkit for mass spectrometry and proteomics. *Nature Biotechnology,*

106 *30,* 918-920.

107 Dunn, W. B., Erban, A., Weber, R. J., Creek, D. J., Brown, M., Breitling, R., et al. (2013). Mass

108 appeal: metabolite identification in mass spectrometry-focused untargeted metabolomics.

109 *Metabolomics, 9,* 44-66.

110 Haug, K., Salek, R. M., Conesa, P., Hastings, J., de Matos, P., Rijnbeek, M., et al. (2013).

111 MetaboLights—an open-access general-purpose repository for metabolomics studies and

112 associated meta-data. *Nucleic Acids Research, 41,* D781-D786.

113 Kind, T., & Fiehn, O. (2007). Seven Golden Rules for heuristic filtering of molecular formulas

114 obtained by accurate mass spectrometry. *BMC Bioinformatics, 8,* 105.

115   Kind, T., Tsugawa, H., Cajka, T., Ma, Y., Lai, Z., Mehta, S. S., et al. (2017). Identification of small

116      molecules using accurate mass MS/MS search. *Mass Spectrometry Reviews*.

117   Patti, G. J., Yanes, O., & Siuzdak, G. (2012). Innovation: Metabolomics: the apogee of the omics

118      trilogy. *Nature Reviews Molecular Cell Biology, 13*, 263-269.

119   Pedrioli, P. G., Eng, J. K., Hubley, R., Vogelzang, M., Deutsch, E. W., Raught, B., et al. (2004). A

120      common open representation of mass spectrometry data and its application to proteomics

121      research. *Nature Biotechnology, 22*, 1459-1466.

122   Perez-Riverol, Y., Bai, M., da Veiga Leprevost, F., Squizzato, S., Park, Y. M., Haug, K., et al.

123      (2017). Discovering and linking public omics data sets using the Omics Discovery Index. *Nature*

124      *Biotechnology, 35*(5), 406-409.

125   Scheubert, K., Hufsky, F., & Böcker, S. (2013). Computational mass spectrometry for small

126      molecules. *Journal of Cheminformatics, 5*, 12.

127   Sévin, D. C., Kuehne, A., Zamboni, N., & Sauer, U. (2015). Biological insights through nontargeted

128      metabolomics. *Current Opinion in Biotechnology, 34*, 1-8.

129   Sud, M., Fahy, E., Cotter., D., Azam., K., Vadivelu., I., Burant, C., et al. (2016). Metabolomics

130      Workbench: An international repository for metabolomics data and metadata, metabolite

131      standards, protocols, tutorials and training, and analysis tools. *Nucleic Acids Research, 44*(D1),

132      D463-D470.

133   Weber, R. J., Lawson, T. N., Salek, R. M., Ebbels, T. M., Glen, R. C., Goodacre, R., et al. (2017).

134      Computational tools and workflows in metabolomics: An international survey highlights the

135      opportunity for harmonisation through Galaxy. *Metabolomics, 13*(2), 12.

136 **Tables**

137 **Table 1 - Untargeted metabolomics studies selected for reanalysis**

| Study Identifier | ST000403 | ST000326 | ST000220 | MTBLS214 | ST000259 |
|---|---|---|---|---|---|
| Instrument | Thermo Scientific Q-Exactive Orbitrap | Agilent 6530 QTOF | Waters Synapt-G2 Si | AB Sciex TripleTOF 5600 | Bruker MicrOTOF II |
| Sample Layout | 6 groups of 3 replicates[a] | 19 individual samples | 3 groups of 7 replicates | 3 groups of 4-5 replicates | 14 groups of 5-6 replicates |
| Compounds /features (+) | 590 | 962 | 1259 | 18 | 857 |
| 50+ Omissions | + | + | + | + | + |
| Mistakes | + | + | + | + | + |
| Study URL | http://www.metabolomicsworkbench.org/data/DRCCMetadata.php?Mode=Study&StudyID=ST000403 | http://www.metabolomicsworkbench.org/data/DRCCMetadata.php?Mode=Study&StudyID=ST000326 | http://www.metabolomicsworkbench.org/data/DRCCMetadata.php?Mode=Study&StudyID=ST000220 | http://www.ebi.ac.uk/metabolights/MTBLS214 | http://www.metabolomicsworkbench.org/data/DRCCMetadata.php?Mode=Study&StudyID=ST000259 |

138 [a]Four of the six groups contained added stable isotope labels. Peaks corresponding to clear stable isotope labeling
139 signals or peaks absent from the two control groups without stable isotope labeling were not considered as possible
140 compound omissions.

141 **Figure Legends**

142 **Figure 1 - Examples of omissions and mistakes in results from study ST000403. (a)**

143 Visualization of a part of one of the raw datafiles. Gray labels correspond to annotations from

144 original results accompanying the study data. Magenta labels correspond to omissions or mistakes.

145 (**b-d**) Mass spectra and extracted ion chromatograms for examples of omissions. (**e**) A mass

146 spectrum and extracted ion chromatograms for examples of mistakes. An in-source fragment ion

147 and an isotopic peak of a multimer of HEPES were incorrectly identified as different compounds.

148 Peaks of related ions for a given compound in plots of mass spectra are highlighted in magenta.

149 Color coding for groups of replicate samples in extracted ion chromatograms is the same as in

150 Supplementary Figure 6.