

# Automated Assembly of Species Metabolomes through Data Submission into a Public Repository

Reza M Salek<sup>1</sup>, Pablo Conesa<sup>1</sup>, Keeva Cochrane<sup>1</sup>, Kenneth Haug<sup>1</sup>, Mark Williams<sup>1</sup>,  
Namrata Kale<sup>1</sup>, Pablo Moreno<sup>1</sup>, Kalai Jayaseelan<sup>1</sup>, Jose Ramon Macias<sup>1</sup>, Venkata  
Chandrasekhar<sup>1</sup> Nainala, Robert D. Hall<sup>4</sup>, Laura K. Reed<sup>2</sup>, Mark R. Viant<sup>3</sup>, Claire  
O'Donovan<sup>1</sup> and Christoph Steinbeck<sup>1,5,\*</sup>

<sup>1</sup> European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-  
EBI), Wellcome Genome Campus, Hinxton, Cambridge CB10 1SD, UK

<sup>2</sup> Department of Biological Sciences, University of Alabama, Tuscaloosa, AL 35487 USA

<sup>3</sup> School of Biosciences, University of Birmingham, Birmingham B15 2TT, UK

<sup>4</sup> Wageningen University & Research, Wageningen Plant Research - Bioscience,  
P.O. Box 16, 6700AA, Wageningen, Netherlands

<sup>5</sup> Institute for Inorganic and Analytical Chemistry, Friedrich-Schiller-University, 07743  
Jena, Germany

\* Correspondence address: Institute for Inorganic and Analytical Chemistry,  
Friedrich-Schiller-University, Lessingstrasse 8, 07743 Jena, Germany. Tel:+49-3641-  
948171; Fax: +49-3641-948172; E-mail: christoph.steinbeck@uni-jena.de

## Abstract

Following similar global efforts to exchange genomic and other biomedical data, global databases in metabolomics have now been established. MetaboLights, the first general purpose, publically available, cross-species, cross-application database in metabolomics, has become the fastest growing data repository at the European Bioinformatics Institute in terms of data volume. Here we present the automated assembly of species metabolomes in MetaboLights, a crucial reference for chemical biology, which is growing through user submissions.

## Keywords

Metabolomics, Databases, Curation, Species Metabolomes

## Background

Following data standardisation efforts in the 1990s and the success of global efforts to exchange genomic [1],[2], proteomic [3], gene expression [4], and other biomedical data, we have now witnessed the emergence of global databases in metabolomics.

In 2012, the European Bioinformatics Institute launched MetaboLights (RRID:SCR\_014663) [5,6], the first general purpose, cross-species, cross-application database in metabolomics, aiming at a similar growth in this remaining large pillar of 'omics sciences [7]. Within the first two years of its inception, MetaboLights became the fastest growing data repository at the EMBL-EBI in terms of data volume. Here we present the automated assembly and growth of species metabolomes in the MetaboLights reference layer, which is largely driven by user

1 submissions. Journals already demand or recommend the deposition of  
2 metabolomics studies in MetaboLights. These include Nature, EMBO, PLOS,  
3  
4 BioMed Central, Frontiers, Metabolomics, and MDPI Metabolites. To the best of our  
5  
6 knowledge, MetaboLights is the only global, general purpose repository that  
7  
8 systematically requires the submission of a metabolites assignment; a requirement  
9  
10 fundamental for the process described here.  
11  
12  
13  
14  
15

## 16 Findings

17  
18  
19  
20  
21 A fundamental, unsolved problem in Metabolomics is the availability of exhaustive  
22  
23 model organism metabolomes. The newly formed Model Organism Metabolomes  
24  
25 task group of the international Metabolomics Society has issued a call to arms to  
26  
27 identify and map all metabolites onto metabolic pathways, and to relate these  
28  
29 pathways across multiple species within the context of evolutionary metabolomics (or  
30  
31 phylometabolomics) [8]. The scale of this endeavour mean the group has prioritized  
32  
33 the deep investigation of established model organism metabolomes in microbial,  
34  
35 plant, and animal biology, promising an avalanche of new metabolic data.  
36  
37 Exponential growth is observed in biological databases, and MetaboLights is no  
38  
39 exception (Fig. 1).  
40  
41  
42  
43  
44  
45  
46

47  
48 Metabolomics datasets submitted to MetaboLights contain lists of metabolites that  
49  
50 have been identified in those respective studies for a given species in a given  
51  
52 biological context. This steady stream of assigned metabolites, together with species  
53  
54 and organism part information, is leading to an evidence-based assembly of  
55  
56 metabolomes for species, with more complete annotations for the model organisms  
57  
58 under investigation worldwide. We believe that this submission-driven assembly,  
59  
60  
61  
62  
63  
64  
65

1 backed by automated and manual quality control, is the only sustainable model for  
2 large-scale species metabolome assembly, and will lead to an indispensable  
3 knowledge base for chemical biology research. This common framework is also  
4 essential to provide the crystallization point to initiate cross-species to cross-division  
5 metabolic analysis of commonality and uniqueness.  
6  
7  
8  
9  
10

11 Studies in MetaboLights are created by researchers in ISA-Tab format, by either  
12 automatically creating datasets from inhouse laboratory information management  
13 systems (rare), or by manually creating ISA-Tab archives with the help of the ISA-  
14 tools suite (common). Naturally, the species coverage of studies follows the  
15 preferences for model species around the globe (Fig. 2).  
16  
17  
18  
19  
20  
21  
22  
23

24 Key to this process is the application of online ontologies from BioPortal, combined  
25 with local controlled vocabularies to ensure correct terms are used to describe  
26 biological samples and experimental factors. Assignment of identified metabolites is  
27 done in Metabolite Identification Files, a bespoke extension to the ISA suite.  
28  
29  
30  
31  
32  
33

34 When the submitter has completed the annotations and the study satisfies all  
35 mandatory validations, the study is flagged as ready for curation [9]. At this stage the  
36 curation team makes any required changes, and the study is ready for review.  
37  
38  
39  
40

41 Journal reviewers are then given a unique URL to access the complete study. When  
42 the journal review process is complete, the study can be made publicly accessible.  
43  
44  
45

46 The traditional model of curated chemical databases scales linearly, both with time  
47 and the number of curators involved in the database assembly.  
48  
49  
50

51 In contrast, the MetaboLights reference layer grows via a two-tiered approach for  
52 assembling metabolome information. We collect historical information about  
53 metabolites found in species from the primary literature, and link it with the  
54 experimental annotations submitted to MetaboLights, thus generating many of the  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1 data points for common and rarer species. In a similar manner as before, this  
2 manually curated data process grows linearly with the number of curators working on  
3  
4  
5 it.

6  
7 The second source for metabolome information is submissions from the community,  
8  
9 triggered by their commitment to provide open access data, and/or the requirement  
10  
11 from funders and publishers to deposit data in an open and accessible manner. The  
12  
13 sustainability and efficiency of this second tier is the key argument of this article.  
14

15  
16 Fig. 3 shows the current distribution of metabolites per species in the MetaboLights  
17  
18 reference layer. Sorted by frequency, this shows a typical long-tail distribution; a  
19  
20 few model species well covered, while only a few metabolites are available for most  
21  
22 species. This data includes both metabolites reported in studies and those manually  
23  
24 curated by MetaboLights and ChEBI curated from the literature.  
25  
26  
27  
28  
29  
30

## 31 Conclusion

32  
33 We have established a model, in which information about metabolites in species  
34  
35 metabolomes grows dynamically through submissions to public archives such as the  
36  
37 MetaboLights database. For the first time, this will automatically provide both the information  
38  
39 about which metabolites are found in which species, and the supporting evidence - the  
40  
41 primary spectroscopic data and supporting meta-data -- in a community driven way. In turn,  
42  
43 this will provide up-to-data knowledge bases for fields such as chemical biology,  
44  
45 metabolomics. and biomedicine.  
46  
47  
48  
49  
50

## 51 List of Abbreviations

52  
53  
54  
55 ChEBI, Chemical Entities of Biological Interest; EMBL-EBI, European Molecular Biology  
56  
57 Laboratory-European Bioinformatics Institute.  
58  
59  
60  
61  
62  
63  
64  
65

# Declarations

## Ethics approval

Not applicable

## Consent for publication

Not applicable

## Availability of data and materials

Data underlying the analysis presented here is available without restrictions at in the MetaboLights database [5] (RRID:SCR\_014663).

## Competing interests

None of the authors have any competing interests

## Funding

The development of MetaboLights was funded by the Biotechnology and Biological Sciences Research Council (BBSRC), <http://dx.doi.org/10.13039/501100000268>, Grant Numbers BB/I000933/1 and BB/L024152/1.

## Authors' contributions

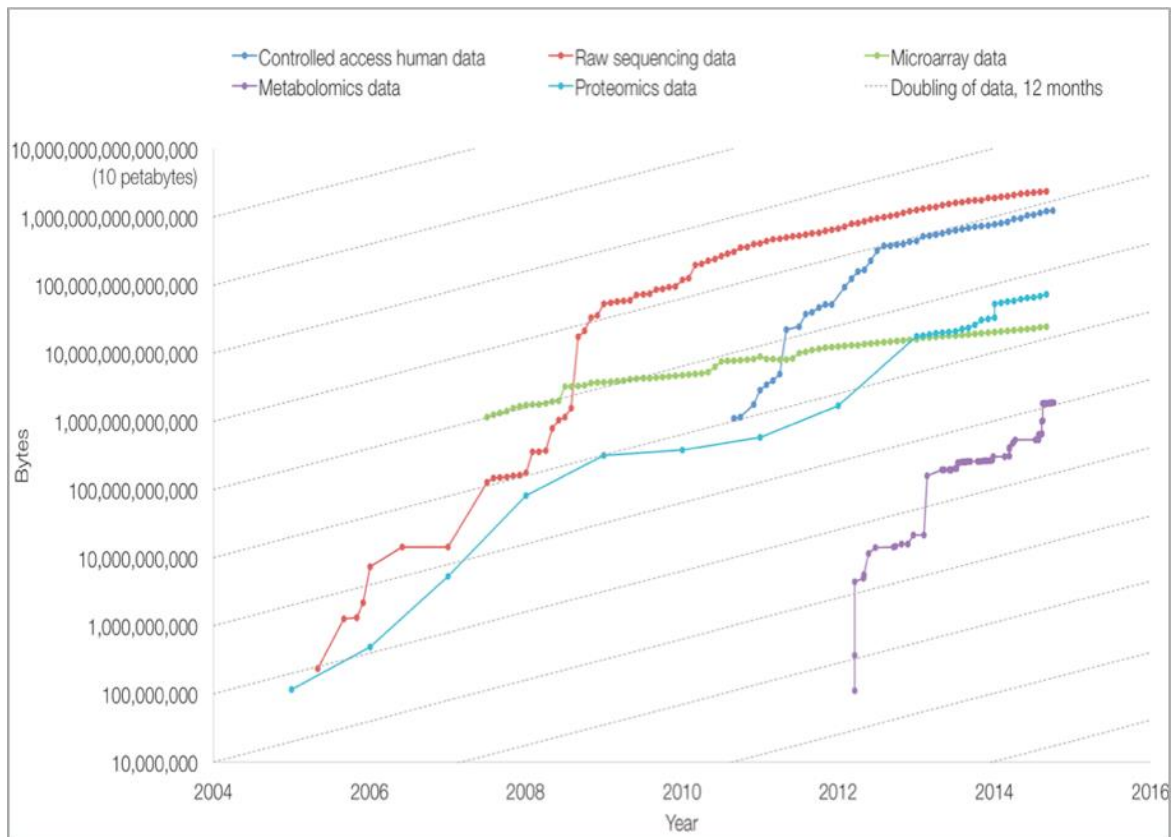
RMS, PC, KC, KH, MW, NK, PM, KJ, JRM, and VCN have developed and curated the MetaboLights database. RDH, LK, MRV, COD, and CS conceived this study and performed the analysis. All authors have read and approved the manuscript.

## Acknowledgements

The authors wish to acknowledge all scientists who deposited data in MetaboLights.

## References

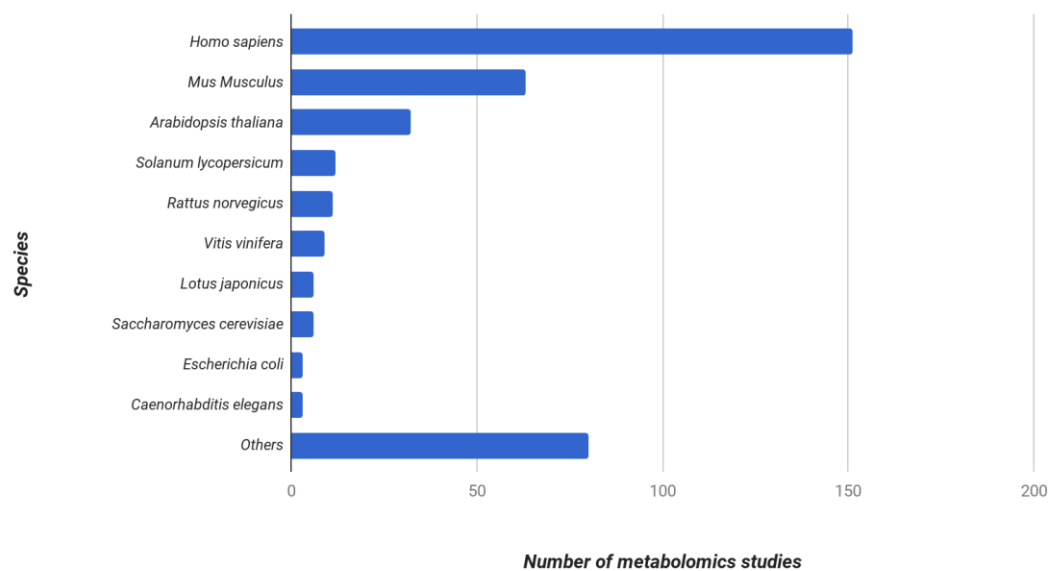
1. Benson DA, Karsch-Mizrachi I, Lipman DJ, Ostell J, Sayers EW. GenBank. *Nucleic Acids Res.* 2009;37:D26–31.
2. Silvester N, Alako B, Amid C, Cerdeño-Tárraga A, Cleland I, Gibson R, et al. Content discovery and retrieval services at the European Nucleotide Archive. *Nucleic Acids Res. Oxford Univ Press*; 2015;43:D23–9.
3. Vizcaíno JA, Deutsch EW, Wang R, Csordas A, Reisinger F, Ríos D, et al. ProteomeXchange provides globally coordinated proteomics data submission and dissemination. *Nat. Biotechnol.* 2014;32:223–6.
4. Kolesnikov N, Hastings E, Keays M, Melnichuk O, Tang YA, Williams E, et al. ArrayExpress update--simplifying data submissions. *Nucleic Acids Res.* 2015;43:D1113–6.
5. MetaboLights. European Molecular Biology Laboratory-European Bioinformatics Institute, Cambridge, UK [Internet]. MetaboLights. [cited 2017 Jul 7]. Available from: <http://www.ebi.ac.uk/metabolights/>
6. Haug K, Salek RM, Conesa P, Hastings J, de Matos P, Rijnbeek M, et al. MetaboLights—an open-access general-purpose repository for metabolomics studies and associated meta-data. *Nucleic Acids Res.* [Internet]. 2012; Available from: <http://nar.oxfordjournals.org/content/early/2012/10/28/nar.gks1004.abstract>
7. Steinbeck C, Conesa P, Haug K, Mahendraker T, Williams M, Maguire E, et al. MetaboLights: towards a new COSMOS of metabolomics data management. *Metabolomics.* Springer; 2012;8:757–60.
8. Edison AS, Hall RD, Junot C, Karp PD, Kurland IJ, Mistrik R, et al. The Time Is Right to Focus on Model Organism Metabolomes. *Metabolites* [Internet]. 2016;6. Available from: <http://dx.doi.org/10.3390/metabo6010008>
9. Salek RM, Haug K, Conesa P, Hastings J, Williams M, Mahendraker T, et al. The MetaboLights repository: curation challenges in metabolomics. *Database* . 2013;2013:bat029.



**Figure 1:** Growth in data repositories at the EMBL-EBI. The graph shows the data volume in each of the repositories over time on a logarithmic scale. Shown are repositories for controlled access human data, raw sequencing data, microarray, proteomics and metabolomics data. Archives were started at different moments in history. Metabolomics shows the steepest growth of all repositories at the EMBL-EBI.

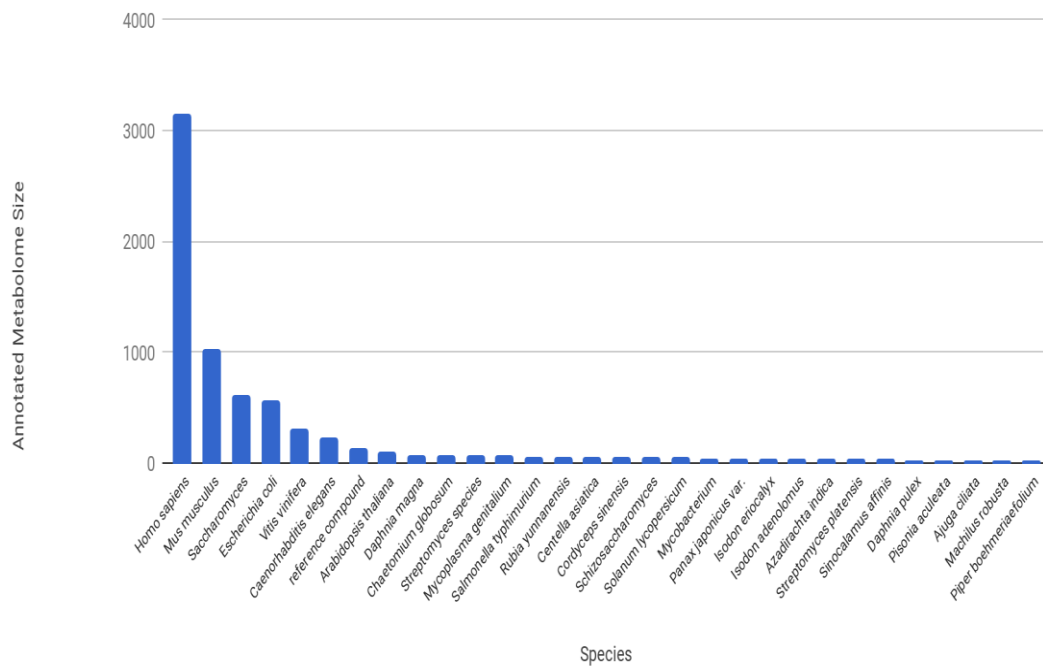


### Number of studies submitted per species

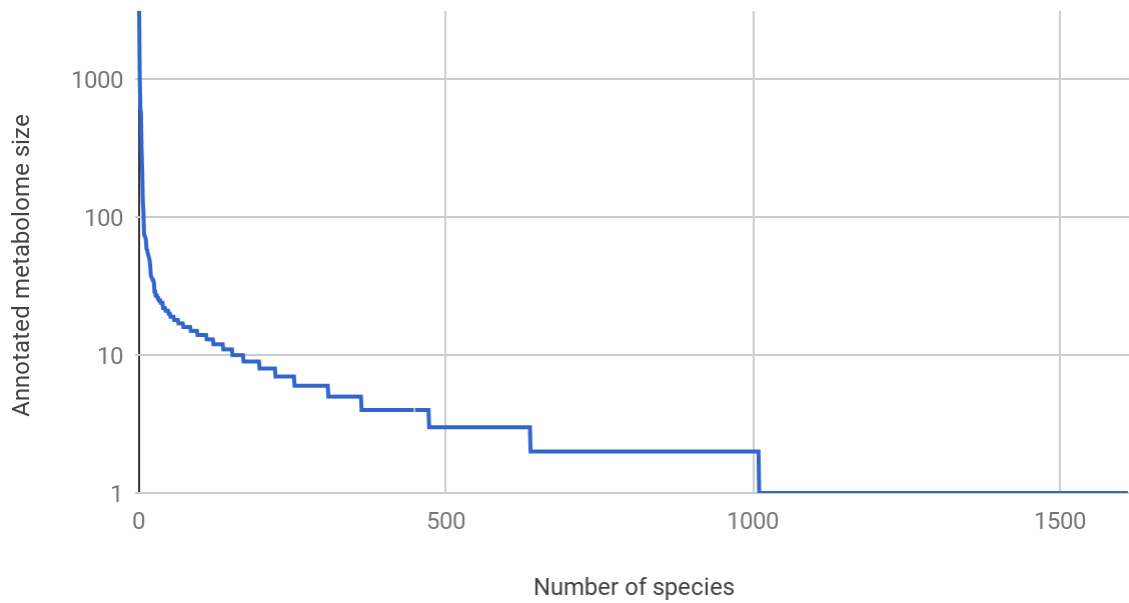


**Figure 2:** Bar chart distribution of Number of studies in MetaboLights by species. The distribution is reflecting the most used model species in biological and biomedical research.

**a**



**b**



**Figure 3:** Long-tail distribution of metabolites per species in the MetaboLights reference layer. A few model species are covered very well, while for the majority of over 1600

species, only a few metabolites were reported. This data covers both metabolites reported in studies as well as those manually added from the literature by MetaboLights and ChEBI curators. a) Truncated version with the 30 most annotated species. b) Full graph.