

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

Reza M Salek¹, Pablo Conesa¹, Keeva Cochrane¹, Kenneth Haug¹, Mark Williams¹, Namrata
5 Kale¹, Pablo Moreno¹, Kalai Jayaseelan¹, Jose Ramon Macias¹, Venkata Chandrasekhar¹
Nainala, Robert D. Hall⁴, Laura K. Reed², Mark R. Viant³, Claire O'Donovan¹ and Christoph
Steinbeck^{1,5,*}

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

10 ² Department of Biological Sciences, University of Alabama, Tuscaloosa, AL 35487 USA

³ School of Biosciences, University of Birmingham, Birmingham B15 2TT UK

⁴ Wageningen University & Research, Wageningen Plant Research - Bioscience, P.O. Box 16,
6700AA, Wageningen, The Netherlands

⁵ Institute for Inorganic and Analytical Chemistry, Friedrich-Schiller-University, 07743 Jena,
15 Germany

ABSTRACT

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

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 (EMBL-EBI) launched the MetaboLights database (RRID:SCR_014663) [5], the first general purpose, cross-species, cross-application database in metabolomics, aiming at a similar growth in this remaining large pillar of omics sciences [6]. Within the first two years after 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 submissions. Following the Bermuda principles[7], which led to the mandatory deposition of data in repositories such as the European Nucleotide Archive [2] or GenBank [1], a number of journals and publishers already demand or recommend the deposition of metabolomics studies in MetaboLights. These include the Nature Journals,

1 including, of course, this journal, the EMBO journal, PLOS journals, MDPI *Metabolites*,
2 BioMed Central, Frontiers journals and *Metabolomics*. To the best of our knowledge,
3
4 MetaboLights is the only global, general purpose repository which systematically requires
5
6 40 the submission of a metabolites assignment - a requirement which is fundamental for the
7
8 process described here.
9

10 Findings

11
12
13
14
15
16
17 One of the fundamental, unsolved problems in Metabolomics is the availability of exhaustive
18
19 model organism metabolomes. The newly formed Model Organism Metabolomes task group
20
21 45 of the international Metabolomics Society has issued a call to arms to identify and map all
22
23 metabolites onto metabolic pathways and to relate these pathways across multiple species
24
25 within the context of evolutionary metabolomics (or phylometabolomics)[8]. Due to the
26
27 scale of this endeavour, the group has prioritised the deep investigation of established
28
29 model organism metabolomes in microbial, plant and animal biology, promising an
30
31 avalanche of new metabolic data. Exponential growth is observed in biological databases
32
33 50 and MetaboLights makes no exception (Figure 1).
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

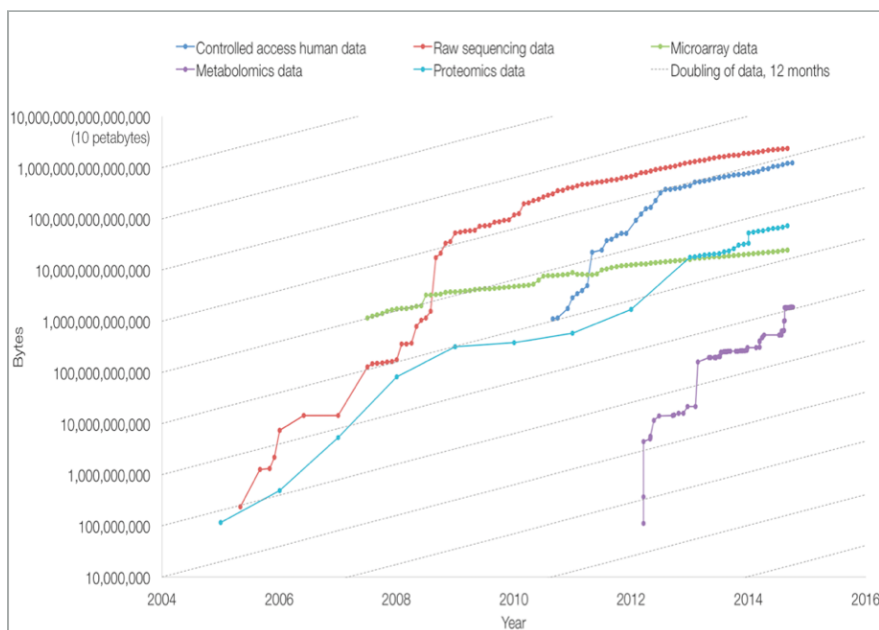


Figure 1: Growth in data repositories at the European Bioinformatics Institute (EMBL-EBI).

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 EMBL-EBI.

Metabolomics datasets submitted to MetaboLights contain lists of metabolites that have been identified in those respective studies for a given species in a given biological context. **This steady stream of assigned metabolites together with species and organism part information is currently leading to an evidence-based assembly of metabolomes for species, with more complete annotations for the model organisms under investigation worldwide.** We believe that this submission driven assembly, backed by automated and manual quality control, will be the only sustainable model for large scale species

1
2 65
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

metabolome assembly and will lead to an indispensable knowledge base for chemical biology research. This common framework is also essential to provide the crystallization point for initiating cross-species to cross-division metabolic analysis of commonality and uniqueness.

Studies in MetaboLights are created by researchers in ISA-TAB format [9], using either automatic creation of datasets from inhouse LIMS systems (rare) or manual creation of ISA-TAB archives with the help of the ISA-tools suite (common). Naturally, the species coverage of studies follows the preferences for model species around the globe (Figure 2).

Key to this process is the application of online ontologies from BioPortal, combined with local controlled vocabularies to ensure correct terms are used to describe biological samples and experimental factors. These controlled vocabularies include the NCBI taxonomy [10] for species identification, the BRENDA tissue ontology [11] and the Experimental Factor Ontology (EFO) [12]. The assignments of identified metabolites is done in Metabolite Identification Files (MAF), a bespoke extension to the ISA suite.

When the submitter has completed the annotations and the study satisfies all mandatory validations, the study is flagged as ready for curation [13]. At this stage the curation team will make any required changes, enabling the study to be ready for review. Journal reviewers are then given a unique URL to access the complete study. After the journal review process completes, the study is ready to be publicly accessible. The traditional model of curated chemical databases scales linearly both with time and the number of curators involved in the database assembly.

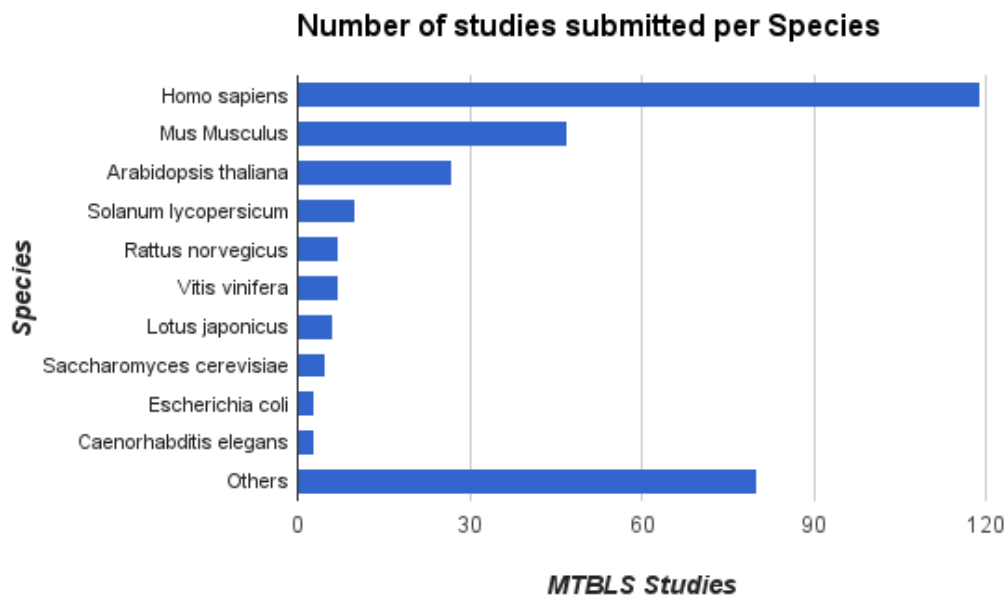


Figure 2: Bar chart distribution of Number of studies in MetaboLights by species. The distribution is reflecting the most used model species in biological and biomedical research.

In contrast, the MetaboLights reference layer grows through a two-tier approach for assembling metabolome information. We collect historical information about metabolites found in species from the primary literature, linking it with the experimental annotations submitted to the MetaboLights, thus generating many of the data points for common and rarer species. In a similar manner as before, this manually curated data process grows linearly with the number of curators working on it.

The second source for metabolome information are submissions from the community triggered by their commitment to open access data and/or the requirement from funders and publishers to deposit data in an open and accessible manner. The sustainability and efficiency of this second tier is the key argument of this article.

Figure 3 shows the current distribution of metabolites per species in the MetaboLights reference layer. Sorted by frequency, this shows a typical long-tail distribution with a few model species being well covered, while for the majority of species, only a few metabolites are available. This data includes both metabolites reported in studies as well as those manually curated by MetaboLights and ChEBI curators from the literature.

1
2
3
4 100
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

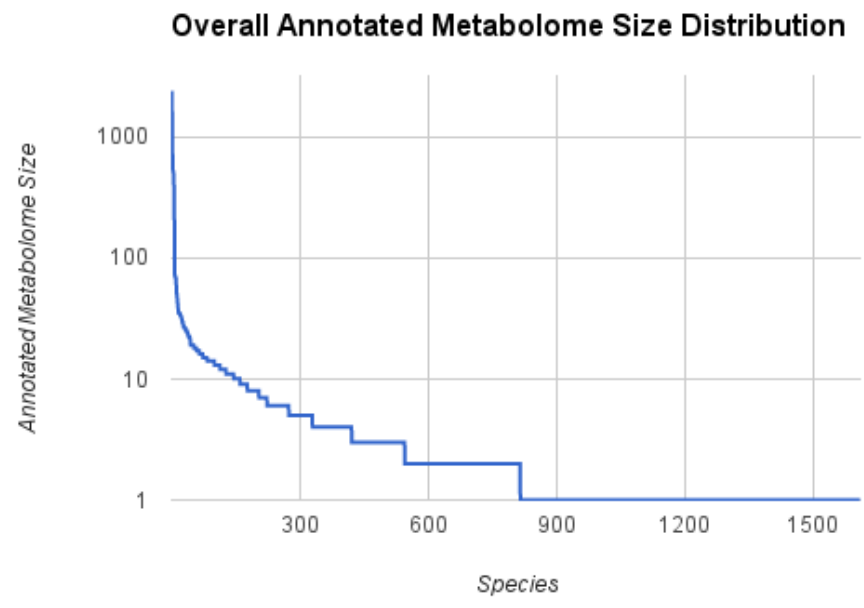
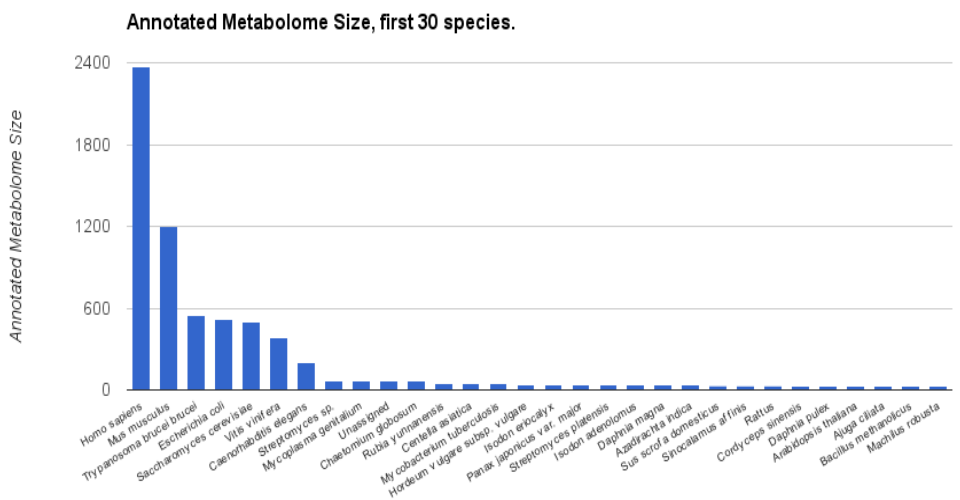


Figure 3: Long-tail distribution of metabolites per species in the MetaboLights reference layer. A few model species are covered very well, while for the majority of over 1600 species, only a few metabolites were reported. This data covers both metabolites reported in studies

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.

Conclusion

In conclusion, we have established a model, where information about metabolites in species metabolomes grows dynamically through submissions to public archives such as the MetaboLights database. For the first time, this will automatically provide both the information about which metabolites are found in which species and also the supporting evidence - the primary spectroscopic data and supporting meta-data - in a community driven way, providing up-to-data knowledge bases for fields such as chemical biology, metabolomics and biomedicine.

Availability of Data and Materials

All data underlying the analysis presented here is available without restrictions at

<http://www.ebi.ac.uk/metabolights> (RRID:SCR_014663).

Competing Interests

None of the authors have any competing interests

1
2
3
4
5
6 **References:**
7
8
9

10 1. Benson DA, Karsch-Mizrachi I, Lipman DJ, Ostell J, Sayers EW. GenBank. *Nucleic Acids Res.*
11
12 125 2009;37:D26–31.
13

14
15 2. Silvester N, Alako B, Amid C, Cerdeño-Tárraga A, Cleland I, Gibson R, et al. Content
16
17 discovery and retrieval services at the European Nucleotide Archive. *Nucleic Acids Res.*
18
19 Oxford Univ Press; 2015;43:D23–9.
20
21

22
23 3. Vizcaíno JA, Deutsch EW, Wang R, Csordas A, Reisinger F, Ríos D, et al. ProteomeXchange
24
25 130 provides globally coordinated proteomics data submission and dissemination. *Nat.*
26
27 *Biotechnol.* 2014;32:223–6.
28
29

30
31 4. Kolesnikov N, Hastings E, Keays M, Melnichuk O, Tang YA, Williams E, et al. ArrayExpress
32
33 update--simplifying data submissions. *Nucleic Acids Res.* 2015;43:D1113–6.
34
35

36
37 5. Haug K, Salek RM, Conesa P, Hastings J, de Matos P, Rijnbeek M, et al. MetaboLights—an
38
39 135 open-access general-purpose repository for metabolomics studies and associated
40
41 meta-data. *Nucleic Acids Res.* [Internet]. 2012; Available from:
42
43 <http://nar.oxfordjournals.org/content/early/2012/10/28/nar.gks1004.abstract>
44
45

46
47 6. Steinbeck C, Conesa P, Haug K, Mahendraker T, Williams M, Maguire E, et al.
48
49 MetaboLights: towards a new COSMOS of metabolomics data management. *Metabolomics.*
50
51 Springer; 2012;8:757–60.
52 140
53
54
55
56

- 1
2 135 7. Schofield PN, Bubela T, Weaver T, Portilla L, Brown SD, Hancock JM, et al. Post-publication
3 sharing of data and tools. *Nature*. 2009;461:171–3.
4
- 5 8. Edison AS, Hall RD, Junot C, Karp PD, Kurland IJ, Mistrik R, et al. The Time Is Right to Focus
6 on Model Organism Metabolomes. *Metabolites* [Internet]. 2016;6. Available from:
7
8 <http://dx.doi.org/10.3390/metabo6010008>
9
- 10 9. Sansone S-A, Rocca-Serra P, Field D, Maguire E, Taylor C, Hofmann O, et al. Toward
11 interoperable bioscience data. *Nat. Genet.* 2012;44:121–6.
12 140
13
14
- 15 10. Federhen S. The NCBI Taxonomy database. *Nucleic Acids Res.* Oxford Univ Press;
16 2012;40:D136–43.
17
18
- 19 11. Gremse M, Chang A, Schomburg I, Grote A, Scheer M, Ebeling C, et al. The BRENDA
20 Tissue Ontology (BTO): the first all-integrating ontology of all organisms for enzyme sources.
21
22
23
24
25
26
27
28
29 145
30
31
32
33
34
35
36
37
38
39
40
41
42
43 145
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65