

# 1 *massDatabase: utilities for the operation of the public compound and pathway database*

2 Xiaotao Shen<sup>1+\*</sup>, Chuchu Wang<sup>2+</sup>, and Michael P. Snyder<sup>1\*</sup>

3 <sup>1</sup>Department of Genetics, Stanford University School of Medicine, Stanford, CA, 94304, USA.

4 <sup>2</sup>Howard Hughes Medical Institute, Stanford University, Stanford, CA 94305, USA.

5

6 <sup>+</sup> These authors contributed equally.

7 <sup>\*</sup> To whom correspondence should be addressed.

8

## 9 **Abstract**

10 **Summary:** One of the major challenges in LC-MS data (metabolome, lipidome, and exposome) is  
11 converting many metabolic feature entries to biological function information, such as metabolite annotation  
12 and pathway enrichment, which are based on the compound and pathway databases. Multiple online  
13 databases have been developed, containing lots of information about compounds and pathways. However,  
14 there is still no tool developed for operating all these databases for biological analysis. Therefore, we  
15 developed *massDatabase*, an R package that operates the online public databases and combines with other  
16 tools for streamlined compound annotation and pathway enrichment analysis. *massDatabase* is a flexible,  
17 simple, and powerful tool that can be installed on all platforms, allowing the users to leverage all the online  
18 public databases for biological function mining. A detailed tutorial and a case study are provided in the  
19 **Supplementary Materials.**

20 **Availability and implementation:** <https://massdatabase.tidymass.org/>.

21 **Contact:** [shenxt@stanford.edu](mailto:shenxt@stanford.edu) and [mepsnyder@stanford.edu](mailto:mepsnyder@stanford.edu)

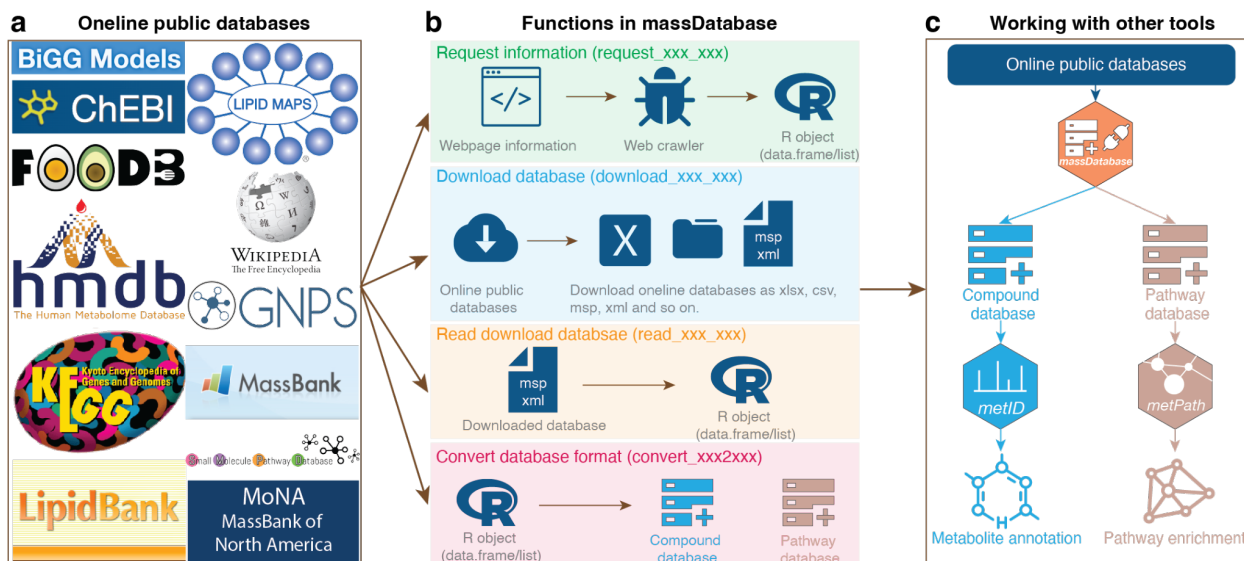
22 **Supplementary information:** Supplementary data are available at *Bioinformatics* online.

23

## 24 **1 Introduction**

25 Liquid chromatography coupled to mass spectrometry (LC-MS) is a comprehensive, unbiased technology  
26 to research small compounds, which has become increasingly popular in dietary, environmental, and  
27 biomedical studies (Wishart, 2016). One of the major challenges in LC-MS data (metabolome, lipidome,  
28 and exposome) is the post-processing of a large number of metabolic feature entries to achieve clear  
29 biological evidence, such as the compound annotation and pathway enrichment. Therefore, the databases  
30 for compounds and pathways are essential for these analyses. Multiple public databases for compounds and  
31 pathways are available online, which benefits the community (Go, 2010). However, the existence of an  
32 automated, multiple compound/pathway query processing package in R is still a demand. So far, although  
33 several R packages have been developed to extract online databases, most of them only support one or  
34 limited databases and have different design concepts and output formats. In addition, they can not be  
35 combined with other existing tools for a straightforward subsequent analysis, which limits their further  
36 applications.

37 Here, we presented the *massDatabase* package to overcome the challenges mentioned above while  
38 accessing the online databases, particularly to (1) support most of the commonly used online public  
39 databases (11 databases, **Table S1**), (2) operate (extracting, downloading, reading, and converting) the  
40 online public databases, and (3) combine the online public databases with existing tools for subsequent  
41 compound annotation, and pathway enrichment analysis (**Fig. 1**).



1  
2 **Fig. 1.** The overview of (a) the online databases that massDatabase support, (b) the functions used to process  
3 databases, and (c) the combination with other tools in the tidyMass project.  
4

## 5 **2 Features and methods**

6 Using *massDatabase*, users can extract compound/pathway information from the online databases (11  
7 databases, **Table S1**) and download them. In addition, *massDatabase* can also be combined with other tools  
8 for metabolite annotation and pathway enrichment analysis. The *massDatabase* can be installed on Mac  
9 OS, Windows, and Linux.  
10

### 11 **2.1 Online databases operation**

12 The functions in *massDatabase* could be grouped into four classes. (1) Request specific information of one  
13 item (compound, pathway, reaction, *etc.*) online using the web crawler, (2) download the corresponding  
14 database, (3) read the downloaded databases (csv, mgf formats, *etc.*) as R object (list or data frame), and  
15 (4) convert the databases to other formats that could be used for other tools (**Fig. 1**).  
16

### 17 **2.2 Combination with other tools**

18 The users can download the online databases and then convert them to the formats supported by the  
19 packages in the tidyMass project using *massDatabase*. Currently, two packages from tidyMass projects  
20 could combine with *massDatabase*. Users can download the compound databases (MS<sup>1</sup> or MS<sup>2</sup> spectra  
21 databases), convert them to the database format in the *metID* package, and then use them for compound  
22 annotation by *metID*. Furthermore, users can also download the pathway databases, convert them to the  
23 pathway database format in the *metPath* package, and then use them for pathway enrichment analysis by  
24 *metPath*.  
25

## 26 **3 Case study**

27 We applied *massDatabase* to a published study from our lab (Liang *et al.*, 2020) as a case study for  
28 exemplifying the value of *massDatabase* in biological function mining by integrating with the online public  
29 databases. The MS<sup>2</sup> spectra databases from HMDB, MassBank, and MoNA were first downloaded and  
30 converted to databases format in *metID*. And the pathway database from KEGG is downloaded and  
31 converted to pathway database format in *metPath*. Then the metabolic feature table was annotated by *metID*

1 which is based on the public databases from *massDatabase* and our in-house library. Then all the annotated  
2 metabolites were used for pathway enrichment analysis using *metPath*. The top enriched pathways include  
3 Steroid hormone biosynthesis, Phenylalanine metabolism, Caffeine metabolism, Linoleic acid metabolism,  
4 Primary bile acid biosynthesis, *etc.*, which are most consistent with the original analysis (**Fig. S1**) (Liang  
5 *et al.*, 2020). These results indicate that *massDatabase* is a powerful tool for utilizing online public  
6 compound and pathway databases for automated and reproducible analysis of LC-MS-based metabolomics  
7 data (**Supplementary Material**).

#### 8 9 **4 Conclusion**

10 *massDatabase* is developed to operate public databases in untargeted LC-MS-based data (metabolome,  
11 lipidome, and exosome). It allows users to extract, download, read databases, and convert database formats  
12 to different formats required by other tools. To our best knowledge, it is the first R package allowing users  
13 to operate most of the commonly used online public databases for subsequent biological function mining.  
14

15 **Funding:** This work received no external funding.

16 **Conflict of Interest:** M.S. is a co-founder and member of the scientific advisory boards of the following:  
17 Personalis, SensOmics, Filtricine, Qbio, January, Mirvie, and Oralome.  
18

#### 19 **References**

- 20 Go,E.P. (2010) Database resources in metabolomics: an overview. *J. Neuroimmune Pharmacol.*, **5**, 18–  
21 30.  
22 Liang,L. *et al.* (2020) Metabolic Dynamics and Prediction of Gestational Age and Time to Delivery in  
23 Pregnant Women. *Cell*, **181**, 1680–1692.e15.  
24 Shen,X. *et al.* (2021) metID: an R package for automatable compound annotation for LC–MS-based data.  
25 *Bioinformatics*, **38**, 568–569.  
26 Shen,X. *et al.* (2022) TidyMass: An Object-oriented Reproducible Analysis Framework for LC-MS Data.  
27 *bioRxiv*, 2022.03.15.484499.  
28 Wishart,D.S. (2016) Emerging applications of metabolomics in drug discovery and precision medicine.  
29 *Nat. Rev. Drug Discov.*, **15**, 473–484.  
30  
31  
32  
33  
34  
35