

TidyMass: An Object-oriented Reproducible Analysis Framework for LC-MS Data

Michael Snyder (✉ mpsnyder@stanford.edu)

Stanford University School of Medicine <https://orcid.org/0000-0003-0784-7987>

Xiaotao Shen

Stanford University School of Medicine <https://orcid.org/0000-0002-9608-9964>

Hong Yan

Yale School of Public Health

Chuchu Wang

Stanford University

Peng Gao

Stanford University <https://orcid.org/0000-0002-4311-584X>

Caroline Johnson

Yale School of Public Health <https://orcid.org/0000-0002-5298-1299>

Brief Communication

Keywords:

Posted Date: March 28th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-1455891/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Nature Communications on July 28th, 2022. See the published version at <https://doi.org/10.1038/s41467-022-32155-w>.

TidyMass: An Object-oriented Reproducible Analysis Framework for LC-MS Data

Xiaotao Shen¹⁺, Hong Yan²⁺, Chuchu Wang³⁺, Peng Gao¹, Caroline H. Johnson^{2*}, and Michael P. Snyder^{1*}

¹ Department of Genetics, Stanford University School of Medicine, Stanford, CA, USA

² Department of Environmental Health Sciences, Yale School of Public Health, New Haven, CT, USA.

³ Howard Hughes Medical Institute, Stanford University, Stanford, CA 94305, USA.

⁺ These authors contributed equally: Xiaotao Shen, Hong Yan, and Chuchu Wang.

^{*} Corresponding authors: Caroline H. Johnson (caroline.johnson@yale.edu) and Michael P. Snyder (mepsnyder@stanford.edu).

Reproducibility and transparency have been longstanding but significant problems for the metabolomics field. Here, we present the tidyMass project (<https://www.tidymass.org/>), a comprehensive computational framework that can achieve the shareable and reproducible workflow needs of data processing and analysis for LC-MS-based untargeted metabolomics. TidyMass was designed based on the following strategies to address the limitations of current tools: 1) Cross-platform utility. TidyMass can be installed on all platforms; 2) Uniformity, shareability, traceability, and reproducibility. A uniform data format has been developed, specifically designed to store and manage processed metabolomics data and processing parameters, making it possible to trace the prior analysis steps and parameters; 3) Flexibility and extensibility. The modular architecture makes tidyMass a highly flexible and extensible tool, so other users can improve it and integrate it with their own pipeline easily.

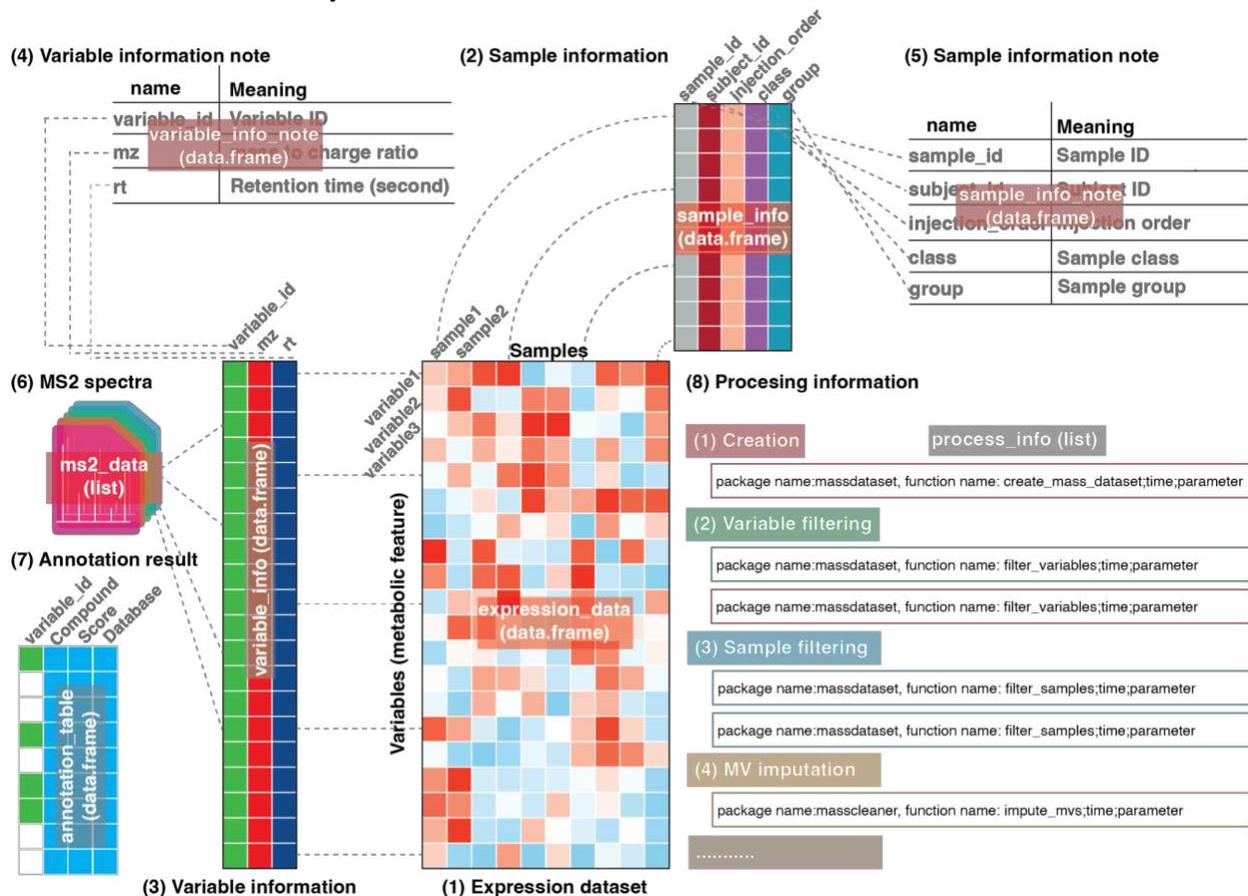
To date, liquid chromatography-mass spectrometry (LC-MS)-based untargeted metabolomics has been proven to be an important tool in environmental, nutrition, and biomedicine research¹. A typical full workflow for LC-MS-based untargeted metabolomics includes sample collection, data acquisition, data analysis, and biological interpretation² (**Fig. S1**). Processing and analyzing high-dimensional metabolomics datasets are challenging, requiring the optimization of multiple steps such as raw data processing, data cleaning, data quality control and assessment, metabolite annotation, statistical analysis, and biological function mining³.

To overcome the challenges of processing and analyzing metabolomics data, the community has developed numerous tools^{4,5}. However, limitations still exist. Commercial tools are expensive and only work on the associated instrument platform, online/GUI tools are user-friendly but cannot take the advantage of the cluster and server computational resources making them impractical for large-scale datasets, open-source tools typically follow limited parts of the whole bioinformatics workflow and have no uniform, specific and traceable format for data input, resulting in a complicated and time-consuming process to prepare data. In addition, different tools with different design concepts and based on different computational platforms make data sharing and reproducible analyses extremely challenging.

Here, we proposed the tidyMass project, an ecosystem of R packages that share an underlying design philosophy, grammar, and data format, which provides a comprehensive, reproducible, and object-oriented computational framework.

We first designed a specific uniform data format (“mass_dataset”) to efficiently store and manage processed untargeted metabolomics data (**Fig. 1**). In the “mass_dataset” class, the expression dataset, metadata of samples and variables are included. Additionally, the datasets in it are automatically synchronous, so when

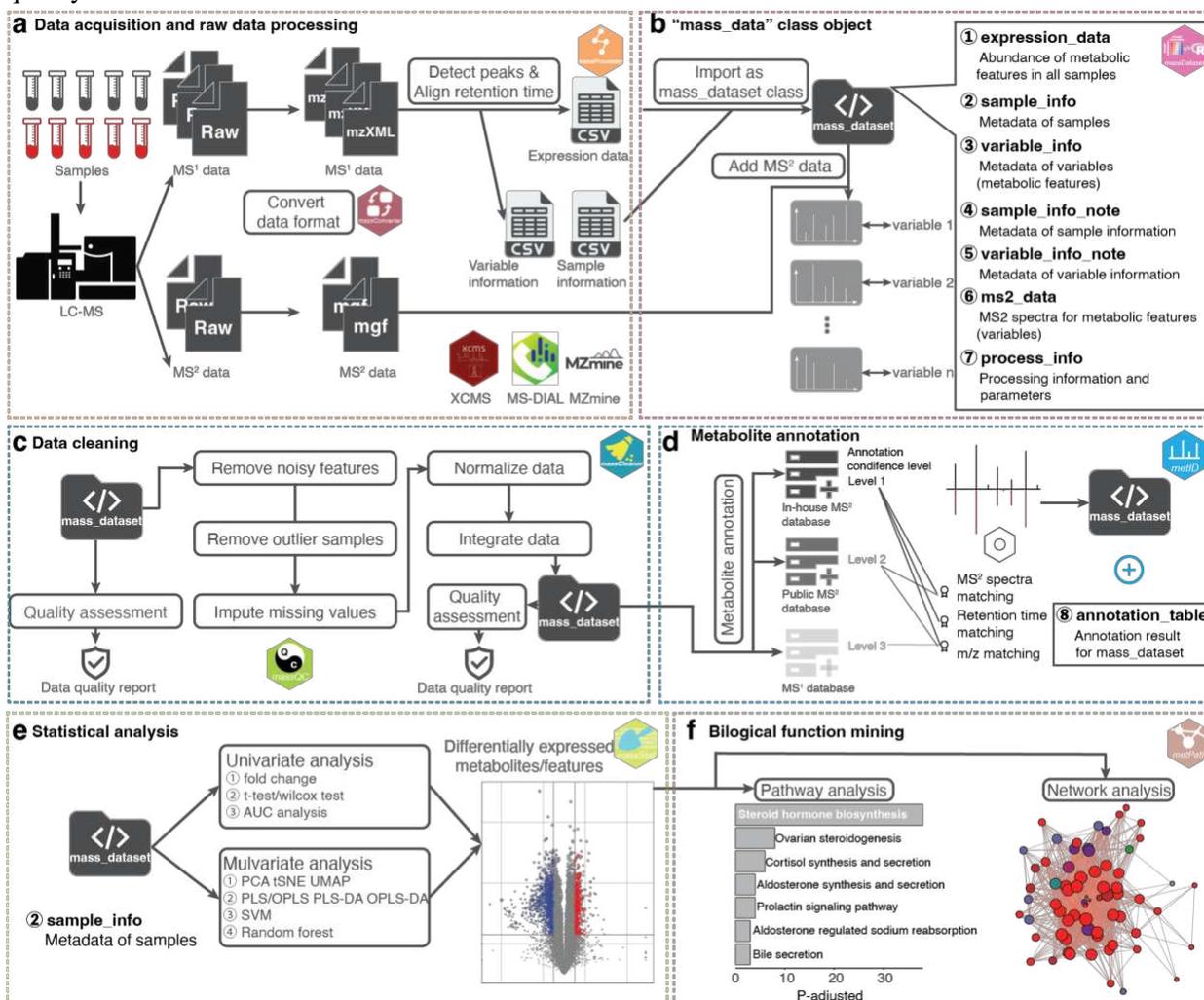
1 the users operate one component, it will automatically propagate the operations across all corresponding
 2 components (**Fig. S2**). This makes it easy to manipulate and maintain the consistency of the data. All the
 3 functions in tidyMass use the “mass_dataset” as their primary input data format, therefore one data format
 4 can be used for all processing and analysis steps (**Fig. 2**). Additionally, the “mass_dataset” class supports
 5 popular tools from other packages, in particular tidyverse, which is one of the most widely used tools for
 6 data science in the R environment⁶ (**Fig. S3**). This design makes the code of tidyMass more universal and
 7 straightforward, which benefits new users as they do not need to adopt new functions. Furthermore, all the
 8 parameters for the processing and analysis are stored in the “mass_data” class object, which makes it
 9 feasible to trace the prior steps and parameters (**Fig. S4**). Briefly, the “mass_dataset” class provides a simple
 10 way to manage and process metabolomics data which sets the foundation for the highly reproducible,
 11 robust, and extendable analytical framework.



12 **Fig. 1 | The “mass_dataset” class and its property.** The “mass_dataset” class is a uniform data format, which is
 13 specifically designed for representing metabolomics data. Most functions in the tidyMass expect this class as their
 14 input format, and all the parameters for the functions can be stored in it.
 15
 16

17 TidyMass provides a set of functions that takes the “mass_dataset” class as the input data format to perform
 18 the whole workflow (**Fig. 2** and **Table S1**). Similar to the concept of tidyverse⁶, tidyMass does not include
 19 all the functions in one package, which is flexible to both users and project managers. TidyMass is a
 20 collection of multiple R packages, where the different packages correspond to different steps of the
 21 workflow (**Fig. 2**). The modular design makes it easy for the user to find appropriate functions, and for
 22 developers to debug and extend it⁷. Briefly, the workflow begins from the package massConverter, which

1 converts MS raw data from different vendors to other formats (**Fig. 2a**). MassConverter depends on the
 2 docker version of msconvert⁸, making it possible to use it on all computational platforms. So the data
 3 conversation can also be integrated with other processing and analysis steps in one code script, which makes
 4 the end-to-end reproducible analysis possible. Next, raw data processing, peak picking and grouping are
 5 performed by the massProcessor package based on XCMS⁹, an object (“mass_dataset” class) is generated
 6 for subsequent analysis in this step. Before moving forward to statistical analysis, data cleaning is
 7 performed to remove unwanted variation by the massCleaner package¹⁰, which carries out noisy features
 8 and outlier sample removal, missing value imputation, data normalization and integration. In the next step,
 9 the metID package performs metabolite annotation using in-house or public databases¹¹. All the statistical
 10 analyses are aimed at finding the potential differentially expressed metabolites using the massStat
 11 package¹². Finally, pathway enrichment analysis is implemented to identify biological functions using the
 12 metPath package. Notably, in any step of the workflow, the massQC package can be used to assess the data
 13 quality.



14 **Fig. 2 | Analysis workflow of tidyMass.**

15 (a) Raw data processing using massConverter and massProcessor. (b) The “mass_dataset” class provides a uniform
 16 data format and object-oriented workflow (massDataset). (c) Data cleaning using massCleaner. (d) Metabolite
 17 annotation using metID. (e) Statistical analysis using massStat. (f) Biological function mining using metPath.
 18
 19

1 Data sharing and reproducible analysis are of utmost importance to avoid biased findings¹³. Unfortunately,
2 reproducibility and transparency for metabolomics within the R environment are less satisfactory than for
3 other types of omics data. Multiple tools offer different parameters, options, and output formats for users.
4 TidyMass is designed to achieve reproducibility and transparency by two aspects. First, the object-oriented
5 class makes it easy to share the data and trace the processing information¹⁴. Second, with the uniform data
6 format and modular design, the users can seamlessly combine all the processing and analyzing steps in an
7 integrative manner in one code script (*e.g.*, Rmarkdown, notebook). In addition, all the steps are optional
8 and the order of execution is customizable, which means that the users can create and optimize customized
9 sharable and reproducible pipelines based on their experimental design and aims. Furthermore, as docker
10 technology is more and more popular in reproducible analysis, we also provide a docker version of
11 tidyMass, containing a R/Rstudio environment and all the tidyMass packages, which makes it possible for
12 users to share all code, data, and even analysis environment based on tidyMass.

13 To demonstrate the application of tidyMass for the processing of metabolomics data, we used data from
14 colorectal cancer (CRC) patient tissues to identify metabolites of CRC by sex of the patient¹⁵ (**Table S3**,
15 **Fig. S5**). First, raw data were converted to mzML format through ProteoWizard⁸, followed by
16 massProcessor to extract the metabolic features. Features with more than 20% missing values (MV) in QC
17 samples or more than 50% MVs in all the study groups were considered as noisy features and were removed.
18 K-nearest neighbors (KNN) was applied for MV imputation, and support vector regression (SVR) enabled
19 data normalization using massCleaner (**Fig. S6**). For metabolite annotation, two in-house databases were
20 constructed using meID, that contain 71 and 55 metabolites in HILIC and RPLC modes, respectively. The
21 databases contain the accurate mass and experimental retentional time of metabolites. A public database¹¹
22 was also used for metabolite annotation. Finally, the redundant annotations were removed based on the
23 annotation score¹¹, and 74 metabolites were identified using the in-house database and up to metabolomics
24 standards initiative (MSI) level 2¹⁶. Only the annotations with level 2 were used for subsequent analysis.
25 We then detected the differentially expressed metabolites between tumor tissues compared to normal
26 controls for males and females separately, using massStat (**Fig. S7**, **Fig. 3a**). Furthermore, metPath was
27 used for pathway enrichment. In addition to our previous findings wherein sex-related differences were
28 observed in methionine, polyamine, pentose phosphate pathways, methionine metabolism and polyamine
29 metabolism¹⁵, we also observed differential enrichment of additional pathways in tumors from female and
30 male patients (**Fig. 3**). For example, ferroptosis and bile acid synthesis was only enriched in tumors from
31 male patients (**Fig. 3 b**). Glutathione metabolism, the cAMP signaling pathway, cGMP-PKG signaling
32 pathway were all enriched in tumors from female patients, but not from males. In addition, tidyMass
33 expedited the analytical workflow, making it more straightforward to analyze and reproducible using a code
34 script (**Supplementary Data 1 and 2**). Additionally, a docker image containing the data, code and analysis
35 environment is also provided for more straightforward reproducible analysis
36 (<https://hub.docker.com/r/jaspershen/tidymass-case-study>).

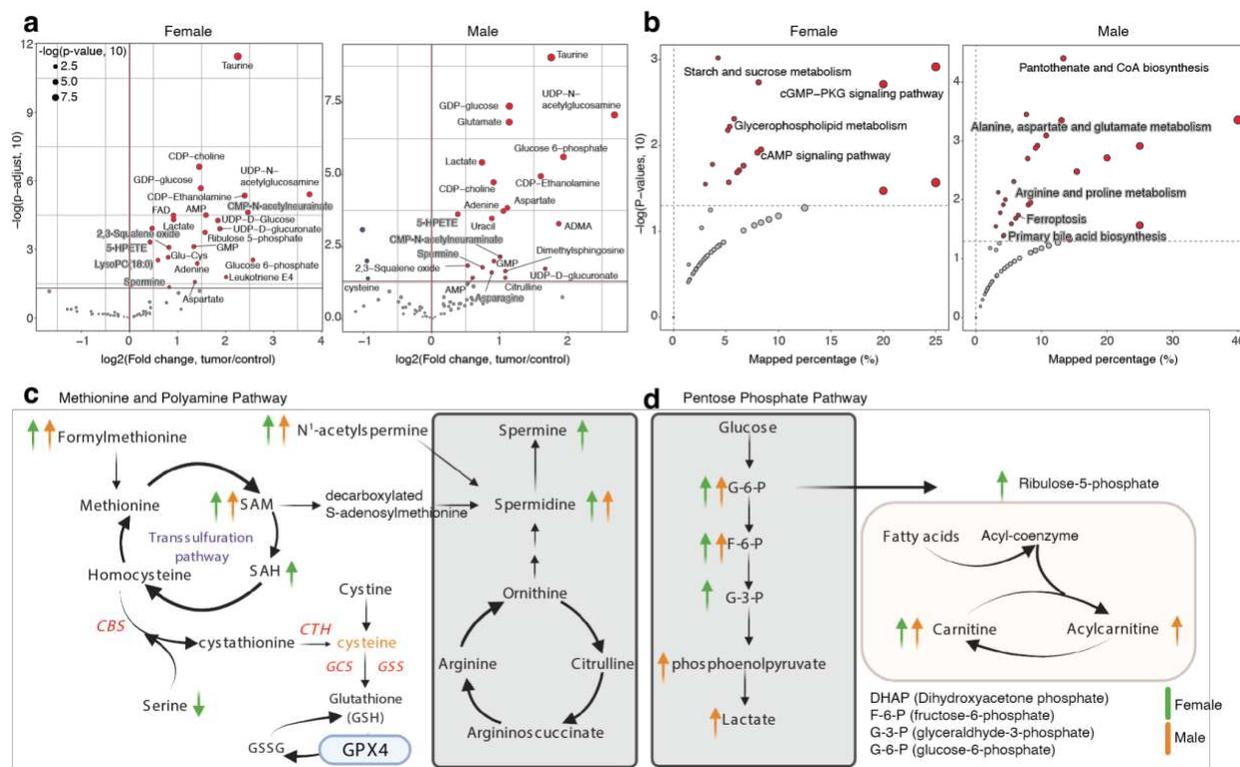


Fig. 3 | Biological function mining for the case study. (a) Volcano plots to show the differentially expressed metabolites. (b) Pathway enrichment analysis. Only the sex difference pathways were labeled. Sex differences in (c) methionine and polyamine pathways, and (d) pentose phosphate pathway metabolism.

In summary, the tidyMass project is an ecosystem of R packages that share an underlying design philosophy, grammar, and uniform data format, which provides a comprehensive, transparent, reproducible, and object-oriented computational framework for LC-MS-based metabolomics data processing and analysis within the R environment. As such, a complete website for tidyMass is publicly available (<https://www.tidymass.org/>). TidyMass can provide great benefit for the metabolomics field, particularly in the two following aspects. 1) Data sharing, tracing, and reproducible analyses. TidyMass provides a specific uniform data format and a whole object-oriented workflow, including a docker image, making data sharing, tracing, and reproducible analysis more straightforward, providing metabolomics researchers the ability to share and repeat analysis feasibly. 2) Flexibility and extensibility. The object-oriented and modular design concept allows for the easy integration of other tools with tidyMass, therefore making tidyMass flexible and extensible within the metabolomics community. An example that illustrates the functions of the tidyverse can be located here: https://massdataset.tidymass.org/articles/tidyverse_verse. However, as a fast-growing field, some widely used metabolomics tools are not wrapped up or supported in tidyMass, such as GNPS¹⁷ and MetaboAnalystR¹⁸, it is capable to convert the “mass_dataset” class to the eligible data format for these tools in the future. Meanwhile, as an open-source tool, tidyMass can be easily implemented into the other pipelines.

1 **Methods**

2 **TidyMass project.** TidyMass project is an ecosystem of R packages that share an underlying design
3 philosophy, grammar, and data structure, which utilizes the concept of tidyverse¹⁹. To address the
4 challenges of data sharing, reproducible analysis, and extensibility, we adopted the object-oriented and
5 modular design concepts, which are also leveraged by other tools^{7,14}.

6
7 *Object-oriented workflow.* In tidyMass, the “mass_dataset” class is designed specifically for storing the
8 metabolomics data and relevant metadata. Most of the functions in all the packages use it as the input data
9 and output format. Based on the “mass_dataset” class and the pipeline function (%>%) from tidyverse,
10 tidyMass provides an object-oriented workflow of data processing and analysis, which is clear and more
11 straightforward.

12
13 *Modular design.* In tidyMass, different packages correspond to different steps of the whole workflow for
14 LC-MS-based untargeted metabolomics data processing and analysis. The fundamental functions are placed
15 in one package named massTools, therefore all the other packages can call those functions from it.
16 Additionally, other developers can easily call these functions in their pipeline. Currently, nine packages in
17 total are included to perform the whole workflow, from raw data processing to biological function mining,
18 and the many graphic functions allow users to generate publication-quality graphics
19 (https://massdataset.tidymass.org/articles/ggplot_mass_dataset). Most of the functions and tools that are
20 widely used are included or supported in tidyMass. For functions/tools that are not yet wrapped, it is simple
21 to implement and integrate them with tidyMass. Finally, one package named “tidymass” was developed to
22 easily install and manage all the packages in the project. For each package, a website with function and
23 package-level help documents and reproducible examples was created to guide new users on how to use it.

24
25 *Naming and Coding style.* In tidyMass, we strove to provide concise and meaningful names. To make the
26 tidymass more user-friendly and easier to use, the coding style of tidyMass follows the tidyverse style guide
27 (<https://style.tidyverse.org/>). Briefly, all the names of packages in the tidyMass project start from “mass”
28 or “met” and follow a noun to describe their function. Such as “massCleaner” which is used for data
29 cleaning and “massQC” which is used for data quality assessment. The variable and function names follow
30 the snake case naming policy, using only lowercase letters, numbers, and underscores are used to separate
31 words within a name. Generally, all the variable names are nouns, and function names are verbs.

32
33 *Help document and tutorials.* We provide the function-level, package-level, and pipeline-level help
34 documents and tutorials as a learning guide for tidyMass. For the function-level help document, the users
35 can find it on the “Reference” page on the corresponding website for each package. It is also possible to
36 access them quickly in the R environment using the “?” function. For the package-level and pipeline-level,
37 the websites are created using the “pkgdown” tools for all the packages, the users can find the help
38 document or tutorial on the “Help document” or “Tutorial” page.

39
40 *Do not reinvent the wheel.* When we designed tidyMass, another important rule was that we did not want
41 to create redundant tools which have similar functions with existing tools. For example, when we want to
42 remove variables from “mass_dataset”, the “filter()” function from the dplyr package is more efficient and
43 popular in the R community. We do not need to create a new function to process the “removing features”
44 step. So in tidyMass, we made use of the base or popular functions in R to support “mass_dataset” to operate

1 the same functions (https://massdataset.tidymass.org/articles/tidyverse_verse,
2 https://massdataset.tidymass.org/articles/base_function). This design also makes it easier for new users to
3 adopt tidyMass and reduce their study burden, it also means that the tidyMass code is more readable and
4 shareable.

5
6 *Deployment and Installation.* All the packages in the tidyMass project are open-source and can be accessed
7 publicly. In case the internet is not stable for one code hosting platform, we deployed it in three different
8 code hosting platforms, namely GitHub (<https://github.com/tidymass>), GitLab
9 (<https://gitlab.com/dashboard/projects>), and Gitee (<https://gitee.com/jaspershen/dashboard/projects>). Any
10 changes will be updated on the three platforms at the same time, so the users can access and install them
11 from at least one platform in any situation.

12
13 **MassDataset package.** The massDataset package is used to provide a uniform data form/structure for LC-
14 MS-based untargeted metabolomics data, relevant metadata, and the corresponding processing parameters
15 (<https://massdataset.tidymass.org/>). Several packages in R provide the object-oriented class for efficient
16 manipulation of sequencing data^{14,20}, and although a similar concept (XCMS3,
17 <https://github.com/sneumann/xcms>) has been also utilized in the metabolomics field¹⁸, there is still no
18 specific uniform data form for all the processing/analysis workflow for LC-MS-based untargeted
19 metabolomics data. Therefore, the massDataset package, the “mass_dataset” class, was specifically
20 designed to store and manage processed metabolomics data and represents this data as an instance of the
21 main data class. This is a key feature of the tidyMass project, all the subsequent wrapped operation functions
22 use this class as their sole or primary input data form.

23
24 *The “mass_dataset” class.* The “mass_dataset” is an S4 object in the R environment that contains nine
25 components (**Fig. 2**), including 1) expression data (expression_data) is a data frame that represents the
26 abundance of all the metabolic features (peaks) in all samples. Each row is a metabolic feature (peak) and
27 each column is a sample. 2) Sample information (sample_info) is a data frame that represents the metadata
28 of samples. The first column is the sample IDs which should be completely identical to the column names
29 of the expression data. Other columns are the attributes of samples, such as subject ID, sample batch,
30 injection order, etc. 3) Variable information (variable_info) is a data frame that represents the metadata of
31 variables (metabolic features or peaks). The first column is the variable IDs which should be completely
32 identical to the row names of the expression data. Other columns are the attributes of variables, such as m/z,
33 rt and mean intensity, etc. 4) Variable information note (variable_info_note) is a data frame that represents
34 the metadata of variable information. 5) Sample information note (sample_info_note) is a data frame that
35 represents the metadata of sample information. 6) MS² spectra (ms2_data) is a list (“ms2_data” class) that
36 is used to store the MS² spectra for peaks. For each spectrum, the parent ion information, MS² spectrum (a
37 data frame with fragment ion m/z and intensity), and the corresponding peak are stored. 7) Annotation result
38 (annotation_table) is a data frame representing the annotation results for variables. 8) Processing
39 information (process_info) is a list (“tidymass_parameters” class) that is used to store the parameters for
40 each processing/analysis step that has been applied on the “mass_dataset” class.

41
42 *Automatic synchronization of components in the “mass_dataset” class.* The components in “mass_dataset”
43 are relevant. For example, the columns of expression data should completely correspond to the rows of
44 sample information, and the rows of expression data should completely correspond to the rows of variable

1 information. When one component in the “mass_dataset” class is modified, other components which are
2 relevant to the changed component will automatically change to keep the consistency of all the components
3 (Fig. S2). This design makes it easier to modify the datasets and keep them consistent.

4
5 *The addition of MS² data to the “mass_dataset” class.* MS² spectra data is important for LC-MS-based
6 untargeted metabolomics data for metabolite annotation. One MS² spectrum is defined by the spectrum
7 information (parent ion information) and MS² spectrum data frame. The spectrum information records the
8 parent ion m/z, retention time, and other information. The MS² spectrum data frame is a matrix with two
9 columns, fragment m/z, and intensity. In the massDataset package, the MS² spectra data can be added to the
10 “mass_dataset” class using the “mutate_ms2()” function. Briefly, the MS² spectra are extracted from the
11 MS² data files (mgf format), and then for each MS² spectrum, it will be assigned to metabolic features based
12 on m/z and retention time matching²¹. To organize and process the MS² data in the “mass_dataset” class, a
13 class named “ms2_data” is designed.

14
15 *Base operation functions for the “mass_dataset” class.* Base operation functions have been provided in the
16 massDataset package to process the “mass_dataset” class. The functions can be divided into four classes
17 (Fig. S3). 1) The first class of functions is used to extract and output datasets in “mass_dataset”. 2) The
18 second class of functions is used to summarize and explore data. 3) The third class of functions is used to
19 preprocess data. For example, add new information, remove samples/variables. 4) The fourth class
20 functions are used to combine or merge two “mass_dataset” class objects. To reduce the difficulty and cost
21 of learning, for the functions which are widely used in R for the same aims but other objects, we wrapped
22 them in massDataset and made the “mass_dataset” class as their input data form. For example, the “filter()”
23 functions from the tidyverse package are widely used in data science to remove eligible variables, so this
24 function is wrapped and users can filter variables from any components in “mass_dataset”.

25
26 *The “tidymass_parameter” class.* To store the parameters for each step that is applied on the
27 “mass_dataset” class, a “tidymass_parameter” class was designed in massDataset. Briefly, four slots are in
28 the “tidymass_parameter” class, namely package name, function name, processing time, and parameter list.
29 The parameter is stored as a list, whose items are specific settings, and the names are arguments. The
30 “tidymass_parameter” classes for all the processing/analysis steps are stored in the “process_info” slot of
31 the “mass_dataset” class and ordered by processing time. Thus, it is possible and easy for the users to trace
32 the processing and analysis for this object. This is another key design in the tidyMass project, which
33 provides the fundamentals for reproducible analysis.

34
35 **MassConverter package.** The massConverter package is used to convert mass spectrometry raw data to
36 different format data (<https://massconverter.tidymass.org/>). MSconvertGUI is the interactive version of the
37 msconvert tool for converting mass spec data files to various formats, which is widely used in the
38 metabolomics field⁸. It also provides the command line version. However, it is software that can only be
39 installed on Windows OS, so cannot be used by Mac OS and Linux users. To achieve a comprehensive
40 reproducible analysis, it is important to do the data converting and record the parameters in the R
41 environment.

42
43 *Docker version of msconvert.* The team provides a docker version of msconvert (pwiz,
44 <https://hub.docker.com/r/chambm/pwiz-skyline-i-agree-to-the-vendor-licenses>), so the massConverter

1 package can convert mass spectrometry raw data to different formats. The users need to install docker based
2 on the official website (<https://www.docker.com/get-started>). Then they pull the pwiz image by using the
3 “docker_pull_pwiz()” function, which will download the pwiz image from the docker hub, therefore it can
4 be used for converting data.

5
6 *Convert data.* Many parameters are included in the mass spectrometry data conversion. The
7 “create_msconvert_parameter()” function is used to set the converting parameters. The detailed converting
8 parameters and their meanings can be found in **Table S2**. After setting the parameters, the
9 “convert_raw_data()” function is used to convert the raw data to other formats. The massConverter package
10 makes it possible to convert the mass spectrometry data using R, and integrate data converting steps with
11 other data processing and analysis in one code file, making the reproducible analysis of metabolomics data
12 more efficient.

13
14 **MassProcessor package.** The massProcessor package is used for mass spectrometry raw data processing,
15 including peak picking and peak grouping based on the widely used XCMS⁹
16 (<https://massprocessor.tidymass.org/>). We have added some new functions to make the results more
17 interpretable. After the processing, a “mass_dataset” class is generated with simple sample information.
18 Then users can add more information directly to it for subsequent processing and analysis using other
19 packages from the tidyMass project. This makes it smoother and more straightforward to combine raw data
20 processing and other processing/analysis steps. In addition, all the graphics from massProcessor, such as
21 “BPC”, “TIC”, and “retention time correction” are generated using the ggplot2 package, which generates
22 high-quality figures for publication. Another important feature of the massProcessor package is that the
23 users can easily extract and score the EIC of all the features and evaluate the quality of features, so can
24 avoid false-positive findings in the subsequent analysis. The raw data processing is optional in the whole
25 workflow, users can use other software/tools to generate peak tables.

26
27 **MassCleaner package.** The massCleaner package is used to do the data cleaning of metabolomics data
28 (<https://masscleaner.tidymass.org/>). The LC-MS-based untargeted metabolomics data always contain
29 different types of bias arising from sample preparation and data acquisition (e.g., contamination, drift in
30 signal intensity.), this is the reason to perform data cleaning as an essential step, which is used to remove
31 unwanted variations. It can be divided into different steps, and some steps are optional, and the orders can
32 be customized based on the study design and aims.

33
34 *Noisy feature removal.* The noisy feature removal can be used based on different rules, according to the
35 experimental aims and design. The functions in massDataset and other packages make it simple to perform
36 the noisy feature removal. For example, the users can define the noisy features as the metabolic features
37 that have missing values more than in 20% QC samples or in 50% subject samples. So the
38 “mutate_variable_na_freq()” function can be added to variable information and then remove the noisy
39 features using the “filter()” function from the dplyr package.

40 *Outlier samples removal.* Outlier samples are a recurrent problem, especially when analyzing large cohorts.
41 Detecting and removing the outlier samples are critical to avoid false positive and false negative findings
42 in the subsequent analysis. Different methods have been used to define and detect outlier samples in
43 tidyMass²². The first rule is the missing value percentage for each sample¹⁰. If one sample with more than
44 50% features is missing values, it means that there may be issues in the sample preparation or data

1 acquisition, so those samples are labeled as an outlier. Other methods are also included to detect outlier
2 samples²³. In brief, all the biological subject samples are used for PCA analysis, then the samples whose
3 principal component 1 (PC1) are more than 6 standard deviations away from the mean value will be labeled
4 as outlier samples. To make this method more robust, we also calculate the median instead of the mean
5 value, and MAD (median absolute deviation) instead of the SD (standard variation) because they are more
6 robust estimators. The last method is based on distance. Instead of using the infinite distance, Mahalanobis
7 distance is a multivariate distance based on all variables (principal components) at once. We use a robust
8 version of this distance, which is implemented in packages `robust` and “`robustbase`” and that is reexported
9 in “`bigutilsr`”. Once the outliers have been detected by different methods, it is easy for users to remove the
10 samples from the “`mass_dataset`” class according to their study aims using the “`filter()`” function.

11
12 *Missing value imputation.* Missing value imputation should be performed after noisy features and outlier
13 sample removal. In `massCleaner`, four widely used methods are implemented to perform missing value
14 imputation: 1) K-nearest neighbors (KNN)²⁴, 2) Bayesian principal component analysis replacement
15 (BPCA)²⁵, 3) `svdImpute`²⁶, 4) random forest imputation (`missForest`)²⁷, 5) zero values, 6) mean values, 7)
16 median values and 8) minimum values. KNN is recommended to impute missing values and set them as
17 default¹⁰.

18
19 *Data normalization and integration.* Data normalization and integration are important to remove the
20 unwanted analytical variations occurring in intra- and inter-batch measurements and to integrate multiple
21 batches forming an integral data set for subsequent statistical analysis²⁸. In the `metCleaner` package, several
22 methods that are widely used are integrated. The methods can be divided into two different classes. The
23 first class is the sample-wise method, including PQN, median, mean, total intensity normalization¹². Total
24 intensity normalization means that all the variable intensity is divided by the total intensity of all the
25 variables in one sample. This method sets the total sum of signals to a constant value for each sample. The
26 median and mean normalization have the same concept. However, these approaches could be hampered.
27 For instance, in the case of large mass differences between samples that may lead to different variable
28 extraction efficiencies between samples¹². The second class is the QC sample-based data normalization,
29 including SVR²⁹, and LOESS³. QC samples are typically generated by mixing aliquots of each subject
30 sample and are regularly analyzed during an experimental run to monitor the stability of the analytical
31 platform and are particularly useful for identifying batch effects. For QC-based data normalization, they
32 require that the first and last injections should be QC samples. The data integration method is used to
33 integrate multiple batch data. In the `massCleaner` function, the QC median, QC mean, subject means, the
34 subject median for each variable (metabolic feature or peak) can be used as the correction factors to
35 integrate batches²⁹.

36
37 **MassQC package.** The `massQC` package is used to assess the data quality of LC-MS-based untargeted
38 metabolomics (<https://massqc.tidymass.org/>). The data quality of metabolomics is visually assessed by
39 several aspects¹⁰. 1) Missing value distribution across samples and/or variables. If one variable (metabolic
40 feature or peak) has more missing values, it means that this variable may be a noisy feature. The same
41 applies to samples that have lots of missing values, it could signify that they are outlier samples that should
42 be removed. 2) RSD (relative standard deviation) for all variables in QC (quality control) samples. Since
43 the QC samples are similar and injected frequently during the data acquisition, the RSD of variables in QC
44 samples can be utilized to evaluate the stability of LC and mass spectrometry. In biomarker discovery, the

1 cutoff is always set as 30%. 3) Intensity of all variables in samples. For QC samples, the median value of
2 the intensity of all variables should be very close. 4) The correlation of QC samples. 5) PCA score plot.
3 The high-quality data assessed by a PCA should show a tight clustering of QC samples relative to the
4 distribution of non-QC samples. In massQC, the users can use the “mass_dataset” as the argument to get
5 the result for each aspect in any step of the whole workflow. In addition, one function named
6 “massqc_report()” can be used to generate an HTML format report including all the results, which is very
7 convenient (https://massqc.tidymass.org/articles/html_qa_report).

8
9 **MetID package.** The metID package is used to perform metabolite annotation based on in-house and
10 available open-source databases (<https://metid.tidymass.org/>)¹¹. It combines information from all major
11 databases for comprehensive and streamlined compound annotation. MetID is a flexible, simple, and
12 powerful tool allowing the compound annotation process to be fully automatic and reproducible. What
13 should be noted is that metID¹¹ was not originally designed for the tidyMass project, so it doesn't support
14 “mass_dataset”. However, it is simple to integrate with tidyMass, which demonstrates the flexibility and
15 extensibility of tidyMass. To integrate metID with the tidyMass project, a function named
16 “annotate_metabolites_mass_dataset()” has been developed to support the “mass_dataset” class. All the
17 annotation results have been organized as a data frame and assigned to “annotation_table” in the
18 “mass_dataset” class. The annotation parameters (matching parameters, the database used, *etc.*) are also
19 assigned to processing information. The users can access the annotation table in “mass_dataset” by using
20 the “extract_annotation_result()” function.

21
22 **MassStat package.** The massStat package is used to perform common statistical analyses within
23 metabolomics analysis (<https://massstat.tidymass.org/>). The massStat package provides efficient tools for
24 the different steps required within the complete data analytics workflow: scaling, univariate analysis,
25 multiple testing correction, multivariate analysis, candidate biomarkers selection, and correlation network
26 analysis.

27
28 *Scaling.* Scaling is a procedure where each variable is modified by a factor and accounts for the different
29 statistical characteristics of each variable. Without scaling, highly abundant compounds tend to dominate
30 the analysis when variance-dependent techniques such as PCA are used. Now in massStat, three commonly
31 used scaling methods are included. Unit-variance scaling (uv) divides each variable by its standard
32 deviation. Pareto scaling, intermediate between no scaling and uv scaling, divides each variable by the
33 square root of the standard deviation. Range scaling divides each variable by its range in all the samples.

34
35 *Univariate analysis.* Commonly used univariate analysis tools have been implemented in tidyMass.
36 Student's t-test (t.test), and Wilcoxon signed-rank test (wilcox.test). The different multiple testing correction
37 methods from p.adjust are also implemented. The fold change, p values, and adjusted p values are directly
38 added to the variable information in “mass_dataset”.

39
40 *Correlation, distance, and correlation network.* Correlation and distance between samples or variables can
41 be calculated using the “cor_mass_dataset()” and “dist_mass_dataset()” functions. The “margin” argument
42 is provided in both functions which requests the sample or variable correlation/distance matrix. The
43 correlation network is widely used to explore the co-expression and co-regulation metabolites, in massStat,
44 the users can obtain a network data format (from ggraph and tidygraph packages) from the “mass_dataset”

1 class object. Then this object can be used for network analysis and visualization using the powerful network
2 analysis ecosystem, including ggraph, igraph, and tidygraph.

3
4 *Multivariate analysis.* It is possible to perform various multivariate analyses, such as PCA, PLS, PLS-DA.
5 A typical first-pass unsupervised method used in untargeted LC-MS-based metabolomics is PCA. The score
6 scatter plots, where each sample is depicted as a point, reveal how all samples relate to each other.
7 Supervised methods such as PLS, PLS-DA³⁰, and clustering are also provided.

8
9 **MetPath package.** The metPath package enables pathway enrichment analysis for metabolomics
10 (<https://metpath.tidymass.org/>). At present, metPath provides two commonly used metabolic pathways for
11 this analysis, KEGG³¹, and SMPDB³². To organize and manage the pathway database, a class named
12 “pathway_database” was designed in the metPath package, which is used to store and manage the pathway
13 data. Like the “mass_dataset” class, the “pathway_database” class can be operated by the base and tidyverse
14 functions, which makes it easy to process and manage the pathway database. Then the Hypergeometric test
15 or Fisher's exact test is performed for pathway enrichment. Different visualization methods for enriched
16 pathways are also provided based on ggplot2 to generate high-quality graphics.

17
18 **MassTools package.** The massTools package provides useful tiny functions for mass spectrometry data
19 processing and analysis (<https://masstools.tidymass.org/>). It is a supporting and base package for the
20 tidyMass project. Some functions are universal and may be used and called by different packages, so they
21 are placed in the massTools package, therefore other packages can directly call those functions anytime and
22 anywhere. For example, the MS² spectra matching plot can be used in different places, so it is also placed
23 in the massTools package.

24
25 **TidyMass package.** The tidyMass package is designed to organize and manage all the packages in the
26 tidyMass project (<https://tidymass.tidymass.org/>), allowing for easy installation and loading multiple
27 “tidyMass” packages in a single step. In brief, all the other packages in the tidyMass project are set as the
28 dependent packages of it, so the users can install all the packages in the tidyMass project by only installing
29 the tidymass package. When one or more packages are updated in the tidyMass project, then users can
30 easily check and update them using the tidyMass package. In addition, users can load all the packages into
31 the R environment by only loading the tidyMass package.

32
33 **Extend tidyMass project.** An increasing number of data processing and analysis tools are being developed
34 within the field of metabolomics. This could be problematic, as the integration of these functions and tools
35 is needed to enable their use in tidyMass. However, the specific and uniform data form (“mass_dataset”
36 class) simplifies the integration of tools that are not wrapped in tidyMass for developers. In fact, in
37 tidyMass, the R base function, tidyverse, and metID package have been integrated with the “mass_dataset”
38 class. In brief, the function should change the “mass_dataset” as its supporting object, and then call the
39 function to process or analyze. A protocol is available to show how to make a function that supports the
40 “mass_dataset” class (https://massdataset.tidymass.org/articles/based_on_mass_dataset). In addition, it is
41 easy to integrate tidyMass with other pipelines. For example, xcmsrocker is an open-source project
42 (<https://github.com/yufree/xcmsrocker>) which was created and maintained by Dr. Miao Yu, this project
43 houses various R packages for LC-MS-based metabolomics data processing and analysis, and tidyMass
44 was recently implemented into this project. Another example is the Stanford Data Ocean

1 (<https://innovations.stanford.edu/sdo>), which is a cloud-based computation platform for multi-omics data
2 processing and analysis, and tidyMass is also implemented onto it.

3
4 **Data preparation for tidyMass.** TidyMass is a flexible pipeline that utilizes the modular design concept,
5 which means that the user can perform a comprehensive and full data processing workflow for
6 metabolomics or can choose to perform various or multiple steps of the workflow.

7
8 *Data preparation for massProcessor.* If the users use the massProcessor package for raw data processing,
9 the mzXML (or mzML) data format should be prepared. All the mzXML format files should be placed in
10 different folders according to their class or group. For example, QC samples and blank samples should be
11 placed into folders named “QC” and “Blank” folders, respectively. Biological subject samples can be placed
12 in a folder named “Subject” or placed into different folders that are named according to the class of samples,
13 for example, “Control” or “Case”.

14
15 *Data preparation for other packages.* The users can also use other software to perform raw data processing
16 to generate the peak (metabolic feature) table, such as MS-DIAL^{32,33}, mzMine³⁴, etc. Then the data can be
17 prepared and the “create_mass_dataset()” function is used to generate the “mass_dataset” class object.
18 These files are required for the “create_mass_dataset()” class. The first file is “expression_data” which is a
19 matrix to store the abundance for each variable in each sample. The column is a sample, and the row is
20 variable. The second file is “sample_info” which is a matrix to store the metadata of samples. What should
21 be noted is that the first column is sample ID (sample_id) which is completely identical to the column
22 names of expression data. The third file is “variable_info” which is a matrix to store the metadata of
23 variables. The first column is the variable ID (variable_id) which should be completely identical to the row
24 names of expression data. In addition, the second column and third column should be mass-to-charge ratio
25 (m/z) and retention time (rt, the unit is second), respectively, which are specific spectral information for
26 mass spectrometry data.

27
28 **Reproducible analysis using tidyMass.** One of the most important aims of tidyMass is to improve the
29 reproducible analysis of LC-MS-based untargeted metabolomics data. In tidyMass, the “mass_dataset”
30 class and modular design make it easier for data sharing and reproducible analysis for metabolomics data.

31
32 *Data sharing.* We have enabled a straightforward method for tidyMass users to share their processed data.
33 After preparing the datasets, a “mass_dataset” class object can be generated using the massDataset package,
34 and then users can share the “mass_dataset” class object with collaborators without the need to share
35 multiple files, which is the typical way of sharing this type of data. Collaborators can load the shared
36 “mass_dataset” class object in the R environment and then directly and easily process it using tidyMass.
37 The users can also output all the components in the “mass_dataset” class to xlsx or csv format, and share
38 one or several files of their choosing.

39
40 *Reproducible analysis.* We encourage users to share their data (“mass_dataset” class) and tidyMass pipeline
41 with other collaborators or journals using R script or R markdown files. As the data processing and analysis
42 code is written by R (tidyMass pipeline), it is straightforward for collaborators to easily reproduce the
43 analysis and results. The demo data (“mass_dataset” class) and R code (R markdown) for our demo data
44 have been provided on the tidyMass homepage (<https://www.tidymass.org/start/>). The demo data and R

1 script of the case study presented are also downloadable on the homepage
2 (https://www.tidymass.org/start/demo_data/).

3
4 *Docker image of tidyMass.* A docker image of tidyMass named “tidymass” has been deployed on the docker
5 hub (<https://hub.docker.com/r/jaspershen/tidymass>). This docker image was developed based on the rocker
6 image verse (<https://hub.docker.com/r/rocker/verse>), which contains a Rstudio and R environment, and
7 installed most of the widely used data science packages, such as tidyverse. We installed all the packages in
8 tidyMass with associated dependent packages, the demo datasets and code were also implemented. The
9 new docker image was then built named “tidymass”. The docker version of tidyMass can be used for data
10 analysis by downloading it and then opening the website version Rstudio for data analysis. The “tidymass”
11 image can also be used as a base image for users who want to build a new image to share their analysis
12 environment with other collaborators or reviewers to repeat their analysis and results. A protocol on how
13 to use the docker image of tidyMass is provided on the website of tidyMass
14 (https://www.tidymass.org/start/tidymass_docker/).

15
16 **Sample preparation and analytical conditions for the case study.** All the sample preparation and
17 analytical conditions for the case study can be found in our previous publication¹⁵.

18 19 **Data availability**

20 All the demo data for how to use tidyMass can be accessible on the tidyMass website
21 (<https://www.tidymass.org/>). For the case study, mass spectrometry raw converted data (mzML) for the
22 case study in this paper is accessible on MetaboLights with MTBLS1122 (HILIC positive), MTBLS1124
23 (HILIC negative), MTBLS1122 (RPLC positive) and MTBLS1130 (RPLC negative). The MS² data (mgf)
24 and processed data (“mass_dataset” class) from the massProcessor package are available on the tidyMass
25 project website (https://www.tidymass.org/start/case_study/), and the “mass_dataset” objects are provided
26 as **Supplementary Data 1**.

27 28 **Code availability**

29 All the source code of the tidyMass project is deployed on GitHub (<https://github.com/tidymass>), GitLab
30 (<https://gitlab.com/users/jaspershen/projects>), and Gitee (<https://gitee.com/jaspershen/projects>), and are
31 public under the MIT License; and works on Windows, macOS X, and most Linux distributions. The docker
32 image of tidyMass is hosted on the docker hub (<https://hub.docker.com/r/jaspershen/tidymass>). The code
33 of the case study (Rmarkdown format, https://www.tidymass.org/start/case_study/) is provided as
34 **Supplementary Data 2**.

35 36 **Acknowledgments**

37 We thank Dr. Miao Yu for integrating the tidyMass project with the xcmsrocker project and providing
38 advice on the development of the docker image version of tidyMass. We also thank Dr. Amir Bahmani and
39 Kexin Cha for integrating the tidyMass project with the Stanford Data Ocean project. We also thank Dr.
40 Axel Brunger for the advice on the manuscript. C.H.J was supported by the NCI/NIH under Award Number
41 K12CA215110, NIGMS/NIH under Award Number 1RM1GM141649-01, and American Cancer Society
42 Research Scholar Grant 134273-RSG-20-065-01-TBE.

1 **Author contributions**

2 X.S. and M.P.S. conceived the method and supervised its implementation. X.S. developed the methods,
3 packages, and the docker image. X.S. and C.W. built the websites and wrote the help documents and
4 tutorials. H.Y. provided and prepared the case study data, H.Y. and X.S. analyzed the case study data. X.S.,
5 H.Y., and C.W. prepared the figures. X.S, H.Y., C.W., C.H.J, and M.S.P wrote the manuscript, C.H.J,
6 M.S.P, and P.G. improved the manuscript. All authors contributed to the final manuscript.

7
8 **Competing interests**

9 M.P.S. is a co-founder and member of the scientific advisory board of Personalis, Qbio, January,
10 SensOmics, Protos, Mirvie, NiMo, Onza, and Oralome. He is also on the scientific advisory board of
11 Danaher, Genapsys, and Jupiter. Other authors declare no conflict of interests.

12
13 **Additional information**

14 **Correspondence and requests for materials** should be addressed to X.S. or M.P.S.

15
16 **References**

- 17 1. Wishart, D. S. Emerging applications of metabolomics in drug discovery and precision medicine.
18 *Nat. Rev. Drug Discov.* **15**, 473–484 (2016).
- 19 2. Alseekh, S. *et al.* Mass spectrometry-based metabolomics: a guide for annotation, quantification and
20 best reporting practices. *Nat. Methods* **18**, 747–756 (2021).
- 21 3. Dunn, W. B. *et al.* Procedures for large-scale metabolic profiling of serum and plasma using gas
22 chromatography and liquid chromatography coupled to mass spectrometry. *Nature Protocols* vol. 6
23 1060–1083 (2011).
- 24 4. Cambiaghi, A., Ferrario, M. & Masseroli, M. Analysis of metabolomic data: tools, current strategies
25 and future challenges for omics data integration. *Brief. Bioinform.* **18**, 498–510 (2017).
- 26 5. Misra, B. B. New software tools, databases, and resources in metabolomics: updates from 2020.
27 *Metabolomics* **17**, 49 (2021).
- 28 6. Website. Wickham et al., (2019). Welcome to the Tidyverse. Journal of Open Source Software,
29 4(43), 1686, <https://doi.org/10.21105/joss.01686>.
- 30 7. Rainer, J. *et al.* A Modular and Expandable Ecosystem for Metabolomics Data Annotation in R.
31 *Metabolites* **12**, 173 (2022).
- 32 8. Chambers, M. C. *et al.* A cross-platform toolkit for mass spectrometry and proteomics. *Nat.*

- 1 *Biotechnol.* **30**, 918–920 (2012).
- 2 9. Smith, C. A., Want, E. J., O’Maille, G., Abagyan, R. & Siuzdak, G. XCMS: processing mass
3 spectrometry data for metabolite profiling using nonlinear peak alignment, matching, and
4 identification. *Anal. Chem.* **78**, 779–787 (2006).
- 5 10. Shen, X. & Zhu, Z.-J. MetFlow: an interactive and integrated workflow for metabolomics data
6 cleaning and differential metabolite discovery. *Bioinformatics* **35**, 2870–2872 (2019).
- 7 11. Shen, X. *et al.* metID: an R package for automatable compound annotation for LC–MS-based data.
8 *Bioinformatics* vol. 38 568–569 (2022).
- 9 12. Blaise, B. J. *et al.* Statistical analysis in metabolic phenotyping. *Nat. Protoc.* **16**, 4299–4326 (2021).
- 10 13. Wratten, L., Wilm, A. & Göke, J. Reproducible, scalable, and shareable analysis pipelines with
11 bioinformatics workflow managers. *Nat. Methods* **18**, 1161–1168 (2021).
- 12 14. McMurdie, P. J. & Holmes, S. phyloseq: An R Package for Reproducible Interactive Analysis and
13 Graphics of Microbiome Census Data. *PLoS One* **8**, e61217 (2013).
- 14 15. Cai, Y. *et al.* Sex Differences in Colon Cancer Metabolism Reveal A Novel Subphenotype. *Scientific*
15 *Reports* vol. 10 (2020).
- 16 16. Sumner, L. W. *et al.* Proposed minimum reporting standards for chemical analysis Chemical
17 Analysis Working Group (CAWG) Metabolomics Standards Initiative (MSI). *Metabolomics* **3**, 211–
18 221 (2007).
- 19 17. Wang, M. *et al.* Sharing and community curation of mass spectrometry data with Global Natural
20 Products Social Molecular Networking. *Nat. Biotechnol.* **34**, 828–837 (2016).
- 21 18. Pang, Z., Chong, J., Li, S. & Xia, J. MetaboAnalystR 3.0: Toward an Optimized Workflow for
22 Global Metabolomics. *Metabolites* vol. 10 186 (2020).
- 23 19. Wickham, H. *et al.* Welcome to the Tidyverse. *Journal of Open Source Software* **4**, 1686 (2019).
- 24 20. Sarfraz, I., Asif, M. & Campbell, J. D. ExperimentSubset: An R package to manage subsets of
25 Bioconductor Experiment objects. *Bioinformatics* (2021) doi:10.1093/bioinformatics/btab179.
- 26 21. Shen, X. *et al.* Metabolic reaction network-based recursive metabolite annotation for untargeted

- 1 metabolomics. *Nat. Commun.* **10**, 1516 (2019).
- 2 22. Sun, H., Cui, Y., Wang, H., Liu, H. & Wang, T. Comparison of methods for the detection of outliers
3 and associated biomarkers in mislabeled omics data. *BMC Bioinformatics* **21**, 357 (2020).
- 4 23. BreunigMarkus, M., KriegelHans-Peter, NgRaymond, T. & SanderJörg. LOF. *ACM SIGMOD*
5 *Record* (2000) doi:10.1145/335191.335388.
- 6 24. Moorthy, K., Mohamad, M. & Deris, S. A Review on Missing Value Imputation Algorithms for
7 Microarray Gene Expression Data. *Current Bioinformatics* vol. 9 18–22 (2014).
- 8 25. Oba, S. *et al.* A Bayesian missing value estimation method for gene expression profile data.
9 *Bioinformatics* vol. 19 2088–2096 (2003).
- 10 26. Troyanskaya, O. *et al.* Missing value estimation methods for DNA microarrays. *Bioinformatics* vol.
11 17 520–525 (2001).
- 12 27. Stekhoven, D. J. & Buhlmann, P. MissForest--non-parametric missing value imputation for mixed-
13 type data. *Bioinformatics* vol. 28 112–118 (2012).
- 14 28. De Livera, A. M. *et al.* Statistical methods for handling unwanted variation in metabolomics data.
15 *Anal. Chem.* **87**, 3606–3615 (2015).
- 16 29. Shen, X. *et al.* Normalization and integration of large-scale metabolomics data using support vector
17 regression. *Metabolomics* vol. 12 (2016).
- 18 30. Rohart, F., Gautier, B., Singh, A. & Cao, K.-A. L. mixOmics: An R package for ‘omics feature
19 selection and multiple data integration. *PLoS Comput. Biol.* **13**, e1005752 (2017).
- 20 31. Kanehisa, M. & Goto, S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.* **28**,
21 27–30 (2000).
- 22 32. Jewison, T. *et al.* SMPDB 2.0: big improvements to the Small Molecule Pathway Database. *Nucleic*
23 *Acids Res.* **42**, D478–84 (2014).
- 24 33. Tsugawa, H. *et al.* MS-DIAL: data-independent MS/MS deconvolution for comprehensive
25 metabolome analysis. *Nat. Methods* **12**, 523–526 (2015).
- 26 34. Pluskal, T., Castillo, S., Villar-Briones, A. & Oresic, M. MZmine 2: modular framework for

1 processing, visualizing, and analyzing mass spectrometry-based molecular profile data. *BMC*
2 *Bioinformatics* **11**, 395 (2010).

3

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementaryData2.rmd](#)
- [supplementary.pdf](#)
- [SupplementaryData1.zip](#)
- [NMETHBC48101Acodeflat.pdf](#)
- [NMETHBC48101Aepcflat.pdf](#)