

# The MetaRbolomics book

*Jan Stanstrup, Corey D. Broeckling, Rick Helmus, Nils Hoffmann, Ewy Mathé, Thomas Naake, Luca Nicolotti, Kristian Peters, Johannes Rainer, Reza M. Salek, Tobias Schulze, Emma L. Schymanski, Michael A. Stravs, Etienne A Thévenot, Hendrik Treutler, Ralf J. M. Weber, Egon Willighagen, Michael Witting, Steffen Neumann*

## Contents

<b>Preface</b>	<b>1</b>
<b>1 Introduction</b>	<b>2</b>
1.1 Metabolomics data processing and analysis . . . . .	3
1.2 The R package landscape . . . . .	4
1.3 Dependences and connectivity of metabolomics packages . . . . .	5
<b>2 R-packages for metabolomics</b>	<b>6</b>
2.1 Mass spectrometry data handling and (pre-)processing . . . . .	7
2.2 Metabolite identification with MS/MS data . . . . .	12
2.3 NMR data handling and (pre-)processing . . . . .	13
2.4 UV data handling and (pre-)processing . . . . .	14
2.5 Statistical analysis of metabolomics data . . . . .	15
2.6 Handling of molecule structures and chemical structure databases . . . . .	20
2.7 Network analysis and biochemical pathways . . . . .	21
2.8 Multifunctional workflows . . . . .	23
2.9 User interfaces and workflow management systems . . . . .	25
2.10 Metabolomics data sets . . . . .	26
<b>3 Conclusions</b>	<b>27</b>
<b>4 Appendices</b>	<b>27</b>
Appendix 1: The MSP File Format and package support . . . . .	28
Appendix 2: metaRbolomics dependencies network . . . . .	30

## Preface

# 1 Introduction

## 1.1 Metabolomics data processing and analysis

## 1.2 The R package landscape

### 1.3 Dependences and connectivity of metabolomics packages

## 2 R-packages for metabolomics

## 2.1 Mass spectrometry data handling and (pre-)processing

### 2.1.1 Profile mode and centroided data

### 2.1.2 Direct infusion mass spectrometry data

### 2.1.3 Hyphenated MS and non-targeted data

### 2.1.4 Targeted data and alternative representations of data

### 2.1.5 Additional dimensionality

### 2.1.6 Structuring data and metadata

Table 1: R packages for mass spectrometry data handling and (pre-)processing.

Functionalities	Package	Repo
<b>MS data handling</b>		
Parser for common file formats: mzXML, mzData, mzML and netCDF. Usually not used directly by the end user, but provides functions to read raw data for other packages.	<a href="#">mzR</a>	BioC
Infrastructure to manipulate, process and visualise MS and proteomics data, ranging from raw to quantitative and annotated data.	<a href="#">MSnbase</a>	BioC
Export and import of processed metabolomics MS results to and from the mzTab-M for metabolomics data format.	<a href="#">rmzTab-M</a>	GitHub
Converts MRM-MS (.mzML) files to LC-MS style .mzML.	<a href="#">MRMConverter</a>	GitHub
Infrastructure for import, handling, representation and analysis of chromatographic MS data.	<a href="#">Chromatograms</a>	GitHub
Infrastructure for import, handling, representation and analysis of MS spectra.	<a href="#">Spectra</a>	GitHub
<b>Peak picking, grouping and alignment (LC-MS focussed or general)</b>		
Pre-processing and visualization for (LC/GC-)MS data. Includes visualization and simple statistics.	<a href="#">xcms</a>	BioC
Automatic optimization of XCMS parameters based on isotopes.	<a href="#">IPO</a>	BioC
Parameter tuning algorithm for XCMS, MZmine2, and other metabolomics data processing software.	<a href="#">Autotuner</a>	BioC
Pre-processing and visualization for (LC/GC-)MS data. Includes visualization and simple statistics.	<a href="#">yamss</a>	BioC
Peak picking with XCMS and apLCMS, low intensity peak detection via replicate analyses. Multi-parameter feature extraction and data merging, sample quality and feature consistency evaluation. Annotation with METLIN and KEGG.	<a href="#">xMSanalyzer</a>	SF
Pre-processing and alignment of LC-MS data without assuming a parametric peak shape model allowing maximum flexibility. It utilizes the knowledge of known metabolites, as well as robust machine learning.	<a href="#">apLCMS</a>	SF
Peak detection using chromatogram subregion detection, consensus integration bound determination and Accurate missing value integration. Outputs in XCMS-compatible format.	<a href="#">warpgroup</a>	GitHub
Peak picking for (LC/GC-)MS data, improving the detection of low abundance signals via a master map of m/z/RT space before peak detection. Results are XCMS-compatible.	<a href="#">cosmiq</a>	BioC

Table 1: R packages for mass spectrometry data handling and (pre-)processing. *(continued)*

Functionalities	Package	Repo
m/z detection (i.e. peak-picking) for accurate mass data, collecting all data points above an intensity threshold, grouping them by m/z values and estimating representative m/z values for the clusters; extracting EICs.	<a href="#">AMDORAP</a>	SF
(GC/LC)-MS data analysis for environmental science, including raw data processing, analysis of molecular isotope ratios, matrix effects, and short-chain chlorinated paraffins.	<a href="#">enviGCMS</a>	CRAN
Sequential partitioning, clustering and peak detection of centroided LC-MS mass spectrometry data (.mzXML), with Interactive result and raw data plot.	<a href="#">enviPick</a>	CRAN
PeakpickingwithXCMS. Groups chemically related features beforealignmentacross samples. Additional processing after alignment includes feature validation, re-integration and annotation based on custom database.	<a href="#">massFlowR</a>	GitHub
<b>Isotope labeling using MS</b>		
Analysis of untargeted LC/MS data from stable isotope-labeling experiments. Also uses XCMS for feature detection.	<a href="#">geoRge</a>	GitHub
Correction of MS and MS/MS data from stable isotope labeling (any tracer isotope) experiments for natural isotope abundance and tracer impurity. Separate GUI available in IsoCorrectoRGUI.	<a href="#">IsoCorrectoR</a>	BioC
Extension of XCMS that provides support for isotopic labeling. Detection of metabolites that have been enriched with isotopic labeling.	<a href="#">X13CMS</a>	NA
Analysis of isotopic patterns in isotopically-labeled MS data. Estimates the isotopic abundance of the stable isotope (either 2H or 13C) within specified compounds.	<a href="#">IsotopicLabelling</a>	GitHub
Finding the dual (or multiple) isotope labeled analytes using dual labeling of metabolites for metabolome analysis (DLEMMA) approach, described in Liron [42].	<a href="#">Miso</a>	CRAN
<b>Targeted MS</b>		
Peak picking using peak apex intensities for selected masses. Reference library matching, RT/RI conversion plus metabolite identification using multiple correlated masses. Includes GUI.	<a href="#">TargetSearch</a>	BioC
Pre-processing for targeted (SIM) GC-MS data. Guided selection of appropriate fragments for the targets of interest by using an optimization algorithm based on user provided library.	<a href="#">SIMAT</a>	BioC
Deconvolution of MS2 spectra obtained with wide isolation windows.	<a href="#">decoMS2</a>	NA
Deconvolution of SWATH-MS experiments to MRM transitions.	<a href="#">SWATHtoMRM</a>	NA
Automatic analysis of large scale MRM experiments.	<a href="#">MRMAnalyzer</a>	NA
Tailors peak detection for targeted metabolites through iterative user interface. It automatically integrates peak areas for all isotopologues and outputs extracted ion chromatograms (EICs).	<a href="#">AssayR</a>	GitHub
Targeted peak picking and annotation. Includes Shiny GUI.	<a href="#">peakPantheR</a>	GitHub
Toolkit for working with Selective Reaction Monitoring (SRM) MS data and other variants of targeted LC-MS data.	<a href="#">sRm</a>	GitHub
Deconvolution of SWATH-MS data.	<a href="#">DecoMetDIA</a>	GitHub
Targeted peak picking and annotation. All functions through Shiny GUI.	<a href="#">TarMet</a>	GitHub
<b>GC-MS and GC×GC-MS</b>		



Table 1: R packages for mass spectrometry data handling and (pre-)processing. *(continued)*

Functionalities	Package	Repo
Unsupervised data mining on GC-MS. Clustering of mass spectra to detect compound spectra. The output can be searched in NIST and ARISTO [50].	<a href="#">MSeasy</a>	CRAN
Pre-processing for GC/MS, MassBank search, NIST format export.	<a href="#">erah</a>	CRAN
Pre-processing using AMDIS [53, 54] for untargeted GC-MS analysis. Feature grouping across samples, improved quantification, removal of false positives, normalisation via internal standard or biomass; basic statistics.	<a href="#">Metab</a>	BioC
Deconvolution of GC-MS and GC $\times$ GC-MS unit resolution data using orthogonal signal deconvolution (OSD), independent component regression (ICR) and multivariate curve resolution (MCR-ALS).	<a href="#">osd</a>	CRAN
Corrects overloaded signals directly in raw data (from GC-APCI-MS) automatically by using a Gaussian or isotopic-ratio approach.	<a href="#">CorrectOverloadedPeaks</a>	CRAN
Alignment of GC data. Also GC-FID or any single channel data since it works directly on peak lists.	<a href="#">GCalignR</a>	CRAN
GC-MS data processing and compound annotation pipeline. Includes the building, validating, and query of in-house databases.	<a href="#">metaMS</a>	BioC
Peak picking for GC $\times$ GC-MS using bayes factor and mixture probability models.	<a href="#">msPeak</a>	SF
Peak alignment for GC $\times$ GC-MS data with homogeneous peaks based on mixture similarity measures.	<a href="#">mSPA</a>	SF
Peak alignment for GC $\times$ GC-MS data with homogeneous and/or heterogenous peaks based on mixture similarity measures.	<a href="#">SWPA</a>	SF
Chemometrics analysis GC $\times$ GC-MS: baseline correction, smoothing, COW peak alignment, multiway PCA is incorporated.	<a href="#">RGCxGC</a>	CRAN
Retention time and mass spectra similarity threshold-free alignments, seamlessly integrates retention time standards for universally reproducible alignments, performs common ion filtering, and provides compatibility with multiple peak quantification methods.	<a href="#">R2DGC</a>	GitHub
<b>Flow injection / direct infusion analysis</b>		
Pre-processing of data from Flow Injection Analysis (FIA) coupled to High-Resolution Mass Spectrometry (HRMS).	<a href="#">proFIA</a>	BioC
Flow In-jection Electrospray Mass Spectrometry Processing: data processing, classification modelling and variable selection in metabolite fingerprinting	<a href="#">FIEmspro</a>	GitHub
Processing Mass Spectrometry spectrum by using wavelet based algorithm. Can be used for direct infusion experiments.	<a href="#">MassSpecWavelet</a>	BioC
<b>Other</b>		
Filtering of features originating from artifactual interference. Based on the analysis of an extract of E. coli grown in 13C-enriched media.	<a href="#">credential</a>	GitHub
Wrappers for XCMS and CAMERA. Also includes matching to a spectral library and a GUI.	<a href="#">metaMS</a>	BioC
Processing of peaktables from AMDIS, XCMS or ChromaTOF. Functions for plotting also provided.	<a href="#">flagme</a>	BioC
Parametric Time Warping (RT correction) for both DAD and LC-MS.	<a href="#">ptw</a>	CRAN
R wrapper for X!Tandem software for protein identification.	<a href="#">rTANDEM</a>	BioC

Table 1: R packages for mass spectrometry data handling and (pre-)processing. (*continued*)

Functionalities	Package	Repo
Building, validation, and statistical analysis of extended assay libraries for SWATH proteomics data.	<a href="#">SwathXtend</a>	BioC
Split a data set into a set of likely true metabolites and likely measurement artifacts by comparing missing rates of pooled plasma samples and biological samples.	<a href="#">MetProc</a>	CRAN
Quality of LC-MS and direct infusion MS data. Generates a report that contains a comprehensive set of quality control metrics and charts.	<a href="#">qcrms</a>	GitHub

### 2.1.7 Ion species grouping and annotation

Table 2: R packages for ion species grouping, annotation, molecular formula generation and accurate mass lookup.

Functionalities	Package	Repo
<b>Molecular formula and isotope analysis</b>		
Simulation of and decomposition of Isotopic Patterns.	<a href="#">Rdisop</a>	BioC
Calculation of isotope fine patterns. Also adduct calculations and molecular formula parsing. Web version available at <a href="http://www.envipat.eawag.ch">www.envipat.eawag.ch</a> .	<a href="#">enviPat</a>	CRAN
Molecular formula assignment, mass recalibration, signal-to-noise evaluation, and unambiguous formula selections are provided.	<a href="#">MFAssignR</a>	GitHub
Uses GenForm for molecular formula generation on mass accuracy, isotope and/or MS/MS fragments, as well as performing MS/MS subformula annotation.	<a href="#">GenFormR</a>	GitHub
Checking element isotopes, calculating (isotope labelled) exact monoisotopic mass, m/z values, mass accuracy, and inspecting possible contaminant mass peaks, examining possible adducts in ESI and MALDI.	<a href="#">MSbox</a>	CRAN
<b>MS feature grouping</b>		
Grouping of correlated features into pseudo compound spectra using correlation across samples and similarity of peak shape. Annotation of isotopes and adducts. Works as an add-on to XCMS.	<a href="#">CAMERA</a>	BioC
Grouping of features based on similarity between coelution profiles.	<a href="#">CliqueMS</a>	CRAN
Cluster based feature grouping for non-targeted GC or LC-MS data.	<a href="#">RAMClustR</a>	CRAN
Uses dynamic block summarisation to group features belong to the same compound. Correction for peak misalignments and isotopic pattern validation.	<a href="#">MetTailor</a>	SF
Isotope & adduct peak grouping, homologous series detection.	<a href="#">nontarget</a>	CRAN
Bayesian approach for grouping peaks originating from the same compound.	<a href="#">peakANOVA</a>	NA
Combination of data from positive and negative ionization mode finding common molecular entities.	<a href="#">MScombine</a>	CRAN
Grouping of correlated features into pseudo compound spectra using correlation across sample. Annotation of isotopes and adducts. Can work directly with the XCMS output.	<a href="#">Astream</a>	NA
Navigation of high-resolution MS/MS data in a GUI based on mass spectral similarity.	<a href="#">MetCirc</a>	BioC

Table 2: R packages for ion species grouping, annotation, molecular formula generation and accurate mass lookup. (*continued*)

Functionalities	Package	Repo
Deconvolution of MS/MS spectra obtained with wide isolation windows.	<a href="#">decoMS2</a>	NA
<b>Ion/adduct/fragment annotation</b>		
Bayesian probabilistic annotation.	<a href="#">ProbMetab</a>	GitHub
Isotope & adduct peak grouping, unsupervised homologous series detection.	<a href="#">nontarget</a>	CRAN
Automatic interpretation of fragments and adducts in MS spectra.	<a href="#">InterpretMSSpectrum</a>	CRAN
Molecular formula prediction based on fragmentation.		
Automated annotation using MS2 data or databases and retention time. Calculation of spectral and chemical networks.	<a href="#">compMS2Miner</a>	GitHub
Screening, annotation, and putative identification of mass spectral features in lipidomics. Default databases contain ~25,000 compounds.	<a href="#">LOBSTAHS</a>	BioC
Automated annotation of fragments from MS and MS2 and putative identification against simulated library fragments of ~500,000 lipid species across ~60 lipid types.	<a href="#">LipidMatch</a>	GitHub
Annotation of lipid type and acyl groups on independent acquisition-mass spectrometry lipidomics based on fragmentation and intensity rules.	<a href="#">LipidMS</a>	CRAN
Accurate mass and/or retention time and/or collisional cross section matching.	<a href="#">masstrixR</a>	GitHub
Downloads KEGG compounds orthology data and wraps the KEGGREST package to extract gene data.	<a href="#">omu</a>	CRAN
Paired mass distance analysis to find independent peaks in m/z-retention time profiles based on retention time hierarchical cluster analysis and frequency analysis of paired mass distances within retention time groups. Structure directed analysis to find potential relationship among those independent peaks. Shiny GUI included.	<a href="#">pmd</a>	CRAN

## 2.2 Metabolite identification with MS/MS data

### 2.2.1 MS/MS data handling, spectral matching and clustering

Table 3: R packages for MS/MS data.

Functionalities	Package	Repo
<b>MS2 and libraries</b>		
Tools for processing raw data to database ready cleaned spectra with metadata.	<a href="#">RMassBank</a>	BioC
From RT-m/z pairs (or m/z alone) creates MS2 experiment files with non-overlapping subsets of the targets. Bruker, Agilent and Waters supported.	<a href="#">MetShot</a>	GitHub
Creating MS libraries from LC-MS data using XCMS/CAMERA packages. A multi-modular annotation function including X-Rank spectral scoring matches experimental data against the generated MS library.	<a href="#">MatchWeiz</a>	GitHub
Assess precursor contribution to fragment spectrum acquired or anticipated isolation windows using "precursor purity" for both LC-MS(/MS) and DI-MS(/MS) data. Spectral matching against a SQLite database of library spectra.	<a href="#">msPurity</a>	BioC
Automated quantification of metabolites by targeting mass spectral/retention time libraries into full scan-acquired GC-MS chromatograms.	<a href="#">baitmet</a>	CRAN
MS2 spectra similarity and unsupervised statistical methods. Workflow from raw data to visualisations and is interfaceable with XCMS.	<a href="#">CluMSID</a>	BioC
Import of spectra from different file formats such as NIST msp, mgf (mascot generic format), and library (Bruker) to MSnbase objects.	<a href="#">MSnio</a>	GitHub
Multi-purpose mass spectrometry package. Contains many different functions .e.g. isotope pattern calculation, spectrum similarity, chromatogram plotting, reading of msp files and peptide related functions.	<a href="#">OrgMassSpecR</a>	CRAN
Annotation of LC-MS data based on a database of fragments.	<a href="#">MetaboList</a>	CRAN
<b>In silico fragmentation</b>		
In silico fragmentation of candidate structures.	<a href="#">MetFragR</a>	GitHub
SOLUTIONS for High ReSOLUTION Mass Spectrometry including several functions to interact with MetFrag, developed during the SOLUTIONS project ( <a href="http://www.solutions-project.eu">www.solutions-project.eu</a> ).	<a href="#">ReSOLUTION</a>	GitHub
Uses MetFrag and adds substructure prediction using the isotopic pattern. Can be trained on a custom dataset.	<a href="#">CCC</a>	GitHub
Retention time prediction based on compound structuredescriptors. Five different machine learning algorithms are available to build models. Plotting available to explore chemical space and model quality assessment.	<a href="#">Retip</a>	GitHub

### 2.2.2 Reading of spectral databases

## 2.3 NMR data handling and (pre-)processing

Table 4: R packages for NMR data handling and (pre-)processing.

Functionalities	Package	Repo
<b>Data processing and Analysis</b>		
A tool for processing of 1H NMR data including: Apodization, baseline correction, bucketing, Fourier transformation, warping and phase correction. Bruker FID can be directly imported.	<a href="#">PepsNMR</a>	GitHub
Spectra alignment, peak picking based processing, Quantitative analysis and visualizations for 1D NMR.	<a href="#">speaq</a>	CRAN
Interactive environment based on R-Shiny that includes a complete set of tools to process and visualize 1D NMR spectral data. Processing includes baseline correction, ppm calibration, removal of solvents and contaminants and re-alignment of chemical shifts.	<a href="#">NMRProcFlow</a>	Bitbucket
TheMetaboMateR toolbox covers basic processing and statistical analysis steps including; several spectral quality assessment (such as dealing with baseline distortions, water suppression to quality assessment of shimming and line width) as well as pre-processing (referencing, baseline correction, ... ) to multivariate analysis statistics functions.	<a href="#">MetaboMate</a>	GitHub
<b>Data Analysis and Identification</b>		
Analysis of 1D and 2D NMR spectra using a ROIs based approach. Export to MMCD or uploaded to BMRB for identification.	<a href="#">rNMR</a>	NA
Pre-processing and identification in an R-based GUI for 1D NMR.	<a href="#">rDolphin</a>	GitHub
Bayesian automated metabolite analyser for 1D NMR spectra. Deconvolution of NMR spectra and automate metabolite quantification. Also identification based on chemical shift lists.	<a href="#">BATMAN</a>	RF
“ASICS: an automatic method for identification and quantification of metabolites in complex 1D 1H NMR spectra.”	<a href="#">ASICS</a>	BioC
ASICSdata: 1D NMR spectra for ASICS.	<a href="#">ASICSdata</a>	BioC
<b>NMR and integration with Genomics</b>		
MWASTools: an integrated pipeline to perform NMR based metabolome-wide association studies (MWAS). Quality control analysis; MWAS using various models (partial correlations, generalized linear models); visualization of statistical outcomes; metabolite assignment using STOCSY; and biological interpretation of MWAS results.	<a href="#">MWASTools</a>	BioC
An Integrated Suite for Genetic Mapping of Quantitative Variations of 1H NMR-Based Metabolic Profiles. mQTL-NMR provides a complete metabolite quantitative trait locus (mQTL) mapping analysis pipeline for metabolomic data.	<a href="#">mQTL.NMR</a>	BioC
Handles hyperspectral data, i.e. spectra plus further information such as spatial information, time, concentrations, etc. Such data are frequently encountered in Raman, IR, NIR, UV/VIS, NMR, MS, etc.	<a href="#">hyperSpec</a>	CRAN

## 2.4 UV data handling and (pre-)processing

Table 5: R Packages for UV data handling and (pre-)processing.

Functionalities	Package	Repo
<b>DAD</b>		
Multivariate Curve Resolution (Alternating Least Squares) for DAD data.	<a href="#">alsace</a>	GitHub
Parametric Time Warping (RT correction) for both DAD and LC-MS.	<a href="#">ptw</a>	CRAN
Handles hyperspectral data, i.e. spectra plus further information such as spatial information, time, concentrations, etc. Such data are frequently encountered in Raman, IR, NIR, UV/VIS, NMR, MS, etc.	<a href="#">hyperSpec</a>	CRAN
Projection based methods for preprocessing, exploring and analysis of multivariate data.	<a href="#">mdatools</a>	CRAN
Collection of baseline correction algorithms, along with a GUI for optimising baseline algorithm parameters.	<a href="#">baseline</a>	CRAN

## 2.5 Statistical analysis of metabolomics data

Table 6: R packages for statistical analysis of metabolomics data.

Functionalities	Package	Repo
<b>Sample size</b>		
Estimate sample sizes for metabolomics experiments, (NMR and targeted approaches supported).	<a href="#">MetSizeR</a>	CRAN
<b>Normalization</b>		
Cross-contribution robust multiple standard normalization.	<a href="#">crmn</a>	CRAN
Normalization using internal standards.		
Within and between batch correction of LC-MS metabolomics data using either QC samples or all samples.	<a href="#">batchCorr</a>	GitLab
Normalisation for low concentration metabolites. Mixed model with simultaneous estimation of a correlation matrix.	<a href="#">Metabnorm</a>	SF
A collection of data distribution normalization methods.	<a href="#">Normalizer</a>	NA
Functions for drift removal and data normalisation based on: component correction, median fold change, ComBat or common PCA (CPCA).	<a href="#">intCor</a>	NA
Normalisation using a singular value decomposition.	<a href="#">EigenMS</a>	SF
Normalization based on RUV-random [164].	<a href="#">MetNorm</a>	CRAN
SVR-based normalization and integration for large-scale metabolomics data.	<a href="#">MetNormalizer</a>	GitHub
Drift correction using QC samples or all study samples.	<a href="#">BatchCorrMetabolomics</a>	GitHub
Signal and Batch Correction for Mass Spectrometry	<a href="#">SBCMS</a>	GitHub
Multiple fitting models to correct intra- and inter-batch effects.	<a href="#">MetaboQC</a>	CRAN
Collection of functions designed to implement, assess, and choose a suitable normalization method for a given metabolomics study.	<a href="#">NormalizeMets</a>	CRAN
<b>Exploratory Data Analysis</b>		
A large number of methods available for PCA.	<a href="#">pcaMethods</a>	BioC
Chemometric analysis of NMR, IR or Raman spectroscopy data. It includes functions for spectral visualisation, peak alignment, HCA, PCA and model-based clustering.	<a href="#">ChemoSpec</a>	BioC
Joint analysis of MS and MS2 data, where hierarchical cluster analysis is applied to MS2 data to annotate metabolite families and principal component analysis is applied to MS data to discover regulated metabolite families.	<a href="#">MetFamily</a>	GitHub
<b>Univariate hypothesis testing</b>		
Many methods for corrections for multiple testing.	<a href="#">multtest</a>	BioC
Estimate tail area-based false discovery rates (FDR) as well as local false discovery rates (fdr) for a variety of null models (p-values, z-scores, correlation coefficients, t-scores).	<a href="#">fdrtool</a>	CRAN
GUI for statistical analysis using linear mixed models to normalize data and ANOVA to test for treatment effects.	<a href="#">MetabR</a>	RF
Derives stable estimates of the metabolome-wide significance level within a univariate approach based on a permutation procedure which effectively controls the maximum overall type I error rate at the $\alpha$ level.	<a href="#">MWSL</a>	GitHub
<b>Multivariate modeling and feature selection</b>		

Table 6: R packages for statistical analysis of metabolomics data. *(continued)*

Functionalities	Package	Repo
Various multivariate methods to analyze metabolomics datasets. Main methods include PCA, Partial Least Squares regression (PLS), and extensions to the PLS like PLS Discriminant Analysis PLS-DA and the orthogonal variants OPLS(-DA).	<a href="#">ropls</a>	BioC
Package for performing Partial Least Squares regression (PLS). PPCA, PPCCA, MPPCA.	<a href="#">pls</a>	CRAN
General framework for building regression and classification models.	<a href="#">MetabolAnalyze</a>	CRAN
ASCA, figure of merit (FoM), PCA, Goeman’s global test for metabolomic pathways (Q-stat), Penalized Jacobian method (for calculating network connections), Time-lagged correlation method and zero slopes method. It also includes centering and scaling functions.	<a href="#">caret</a>	CRAN
RF for the construction, optimization and validation of classification models with the aim of identifying biomarkers. Also normalization, scaling, PCA, MDS.	<a href="#">MetStaT</a>	CRAN
PLS-DA, RF, SVM, GBM, GLMNET, PAM.	<a href="#">RFmarkerDetector</a>	CRAN
Recursive feature elimination approach that selects features which significantly contribute to the performance of PLS-DA, Random Forest or SVM classifiers.	<a href="#">OmicsMarkeR</a>	BioC
Find Biomarkers in two class discrimination problems with variable selection methods provided for several classification methods (lasso/elastic net, PC-LDA, PLS-DA, and several t-tests).	<a href="#">biosigner</a>	BioC
Unsupervised feature extraction specifically designed for analysing noisy and high-dimensional datasets.	<a href="#">BioMark</a>	CRAN
Non-parametric method for identifying differentially expressed features based on the estimated percentage of false predictions.	<a href="#">KODAMA</a>	CRAN
Fits multi-way component models via alternating least squares algorithms with optional constraints: orthogonal, non-negative, unimodal, monotonic, periodic, smooth, or structure. Fit models include InDScal, PARAFAC, PARAFAC2, SCA, Tucker.	<a href="#">RankProd</a>	CRAN
Decompose a tensor of any order, as a generalisation of SVD also supporting non-identity metrics and penalisations. 2-way SVD is also available. Also includes PCAn (Tucker-n) and PARAFAC/CANDECOMP.	<a href="#">multiway</a>	CRAN
Fits multi-way component models via alternating least squares algorithms with optional constraints. Fit models include Individual Differences Scaling, Multiway Covariates Regression, PARAFAC (1 and 2), SCA, and Tucker Factor Analysis.	<a href="#">PTAk</a>	CRAN
Performs variable selection in a multivariate linear model by estimating the covariance matrix of the residuals then use it to remove the dependence that may exist among the responses and eventually performs variable selection by using the Lasso criterion.	<a href="#">ThreeWay</a>	CRAN
Performs the O2PLS data integration method for two datasets yielding joint and data-specific parts for each dataset.	<a href="#">MultiVarSel</a>	CRAN
Contains ordination methods such as ReDundancy Analysis (RDA), (Canonical or Detrended) Correspondence Analysis (CCA, DCA for binary explanatory variables), (Non-metric) Multi-Dimensional Scaling ((N)MDS) and other univariate and multivariate methods. Originally developed for vegetation ecologists, many functions are also applicable to metabolomics.	<a href="#">OmicsPLS</a>	CRAN
	<a href="#">vegan</a>	CRAN



Table 6: R packages for statistical analysis of metabolomics data. *(continued)*

Functionalities	Package	Repo
Linear and non-linear Discriminant Analysis methods (e.g. LDA), stepwise selection and classification methods useful for feature selection.	<a href="#">klaR</a>	CRAN
Variable selection methods for PLS, including significance multivariate correlation (SMC), selectivity ratio (SR), variable importance in projections (VIP), loading weights (LW), and regression coefficients (RC). It contains also some other modelling methods.	<a href="#">plsVarSel</a>	CRAN
Predictive multivariate modelling using PLS and Random Forest Data. Repeated double cross unbiased validation and variable selection.	<a href="#">MUVR</a>	GitLab
Biomarker validation for predicting survival. Cross validation methods to validate and select biomarkers when the outcome of interest is survival.	<a href="#">MetabolicSurv</a>	CRAN
Pre-treatment, classification, feature selection and correlation analyses of metabolomics data.	<a href="#">metabolysR</a>	GitHub
<b>Omics Data integration</b>		
Multiple co-inertia analysis of omics datasets (MCIA) is a multivariate approach for visualization and integration of multi-omics datasets. The MCIA method is not dependent on feature annotation therefore can extract important features even when there are not present across all datasets.	<a href="#">omicade4</a>	BioC
STATegRa combines information in multiple omics datasets to evaluate the reproducibility among samples and across experimental condition using component analysis (omicsNPC implements the NonParametric Combination) and clustering.	<a href="#">STATegRa</a>	BioC
Statistical framework supporting many different types of multivariate analyses (e.g. PCA, PLS, CCA, PLS-DA, etc.).	<a href="#">mixOmics</a>	CRAN
STatistics in R Using Class Templates - Classes for building statistical workflows using methods, models and validation objects. Provides mechanism for STATO integration.	<a href="#">STRUCT</a>	GitHub
Multi-omics base classes integrable with commonly used R Bioconductor objects for omics data; container that holds omics results.	<a href="#">MultiDataSet</a>	BioC
Identifies analyte-analyte (e.g. gene-metabolite) pairs whose relationship differs by phenotype (e.g. positive correlation in one phenotype, negative or no correlation in another). The software is also accessible as a user-friendly interface at <a href="http://intlim.bmi.osumc.edu">intlim.bmi.osumc.edu</a> .	<a href="#">IntLIM</a>	GitHub
<b>Missing value imputation</b>		
Mixture-model for accounting for data missingness.	<a href="#">metabomxtr</a>	BioC
Kernel-Based Metabolite Differential Analysis provides a kernel-based score test to cluster metabolites between treatment groups, in order to handle missing values.	<a href="#">KMDA</a>	CRAN
Visualization and imputation of missing values. VIM provides methods for the evaluation and visualization of the type and patterns of missing data. The included imputation approaches are kNN, Hot-Deck, iterative robust model-based imputation (IRMI), fast matching/imputation based on categorical variables and regression imputation.	<a href="#">VIM</a>	CRAN
Graphical user interface for VIM.	<a href="#">VIMGUI</a>	CRAN

Table 6: R packages for statistical analysis of metabolomics data. *(continued)*

Functionalities	Package	Repo
kNN based imputation for microarray data.	<a href="#">impute</a>	BioC
Bootstrap based algorithm and diagnostics for fast and robust multiple imputation for cross sectional, time series or combined cross sectional and time series data.	<a href="#">Amelia</a>	CRAN
Algorithms and diagnostics for the univariate imputation of time series data.	<a href="#">imputeTS</a>	CRAN
Methods for the Imputation of incomplete continuous or categorical datasets. missMDA allows missing data imputation using in categorical, continuous or mixed-type datasets using PCA, CA, a multiple correspondence analysis (MCA) model, a multiple factor analysis (MFA) model or factorial analysis for mixed data (FAMD).	<a href="#">missMDA</a>	CRAN
Random forest based missing data imputation for mixed-type, nonparametric data. An out-of-bag (OOB) error estimate is used for model optimization.	<a href="#">missForest</a>	CRAN
Multivariate imputation by chained equations using fully conditional specifications (FCS) for categorical, continuous and binary datasets. It includes various diagnostic plots for the evaluation of the imputation quality.	<a href="#">mice</a>	CRAN
Missing data imputation using an approximate Bayesian framework. Diagnostic algorithms are included to analyze the models, the assumptions of the imputation algorithm and the multiply imputed datasets.	<a href="#">mi</a>	CRAN
Iterative Gibbs sampler based left-censored missing value imputation.	<a href="#">GSimp</a>	GitHub
<b>Multiple workflow steps</b>		
Missing value imputation, filtering, normalisation and averaging of technical replications.	<a href="#">MSPrep</a>	SF
HCA, Fold change analysis, heat maps, linear models (ordinary and empirical Bayes), PCA and volcano plots. Also log transformation, missing value replacement and methods for normalisation.	<a href="#">metabolomics</a>	CRAN
Cross-contribution compensating multiple internal standard normalisation (ccmn) and remove unwanted variation (ruv2).		
Data processing, normalization, statistical analysis, metabolite set enrichment analysis, metabolic pathway analysis, and biomarker analysis.	<a href="#">MetaboAnalystR</a>	GitHub
Pipeline for metabolomics data pre-processing, with particular focus on data representation using univariate and multivariate statistics. Built on already published functions.	<a href="#">muma</a>	GitHub
Framework for multiomics experiments. Identifies sources of variability in the experiment and performs additional analysis (identification of subgroups, data imputation, outlier detection).	<a href="#">MOFA</a>	BioC
Performs entry-level differential analysis on metabolomics data.	<a href="#">MetaboDiff</a>	GitHub
STRUCT wrappers for filtering, normalisation, missing value imputation, glog transform, HCA, PCA, PLSDA, PLSR, t-test, fold-change, ANOVA, Mixed Effects, post-hoc tests	<a href="#">STRUCTToolbox</a>	GitHub
Data transformation, filtering of feature and/or samples and data normalization. Quality control processing, statistical analysis and visualization of MS data.	<a href="#">pmartR</a>	GitHub
Quality control, signal drift and batch correction, transformation, univariate hypothesis testing.	<a href="#">metabolis</a>	GitHub

Table 6: R packages for statistical analysis of metabolomics data. *(continued)*

Functionalities	Package	Repo
Missing value filtering and imputation, zero value filtering, data normalization, data integration, data quality assessment, univariate statistical analysis, multivariate statistical analysis such as PCA and PLS-D and potential marker selection	<a href="#">MetCleaning</a>	GitHub
Univariate analysis (linear model), PCA, clustered heatmap, and partial correlation network analysis. Based on classes from the Metabase package(Zhu 2019).	<a href="#">ShinyMetabase</a>	GitHub
Outlier detection, PCA, drift correction, visualization, missing value imputation, classification.	<a href="#">MetabolomicsBasics</a>	CRAN
Pre-processing, differential compound identification and grouping, traditional PK parameters calculation, multivariate statistical analysis, correlations, cluster analyses and resulting visualization.	<a href="#">polyPK</a>	CRAN

## 2.6 Handling of molecule structures and chemical structure databases

Table 7: R Packages for molecule structures and chemical structure databases.

Functionalities	Package	Repo
<b>Structure representation and manipulation</b>		
Subset of functions from the Chemistry Development Kit. Provide a computer readable representation of molecular structures and provide functions to import structures from different molecule structure description formats, manipulate structures, visualize structures and calculate properties and molecular fingerprints.	<a href="#">rdck</a>	CRAN
Similar torcdkin functionality and provides more fingerprints and clustering methods and provides additional tools through querying the ChemMine Tools web service.	<a href="#">ChemmineR</a>	BioC
Provides conversion of structure representation through OpenBabel.	<a href="#">ChemmineOB</a>	BioC
Exposes functionalities of the RDKit library, including reading and writing of SF files and calculating a few physicochemical properties.	<a href="#">RRDKit</a>	GitHub
Read and write InChI and InChIKey from and torcdk.	<a href="#">rinchi</a>	GitHub
Maximum Common Substructure Searching using ChemmineR structures.	<a href="#">FmcsR</a>	BioC
Basic cheminformatics functions tailored for mass spectrometry applications, enhancing functionality available in other packages likercdk, enviPat, RMassBank etc.	<a href="#">RChemMass</a>	GitHub
Provides fingerprinting methods forrdck.	<a href="#">fingerprint</a>	CRAN
<b>Database queries</b>		
Calculation of molecular properties.	<a href="#">camb</a>	GitHub
Querying information from PubChem.	<a href="#">Rpubchem</a>	CRAN
Querying information from various web services (CACTUS, CTS, PubChem, ChemSpider) as part of compound list generation.	<a href="#">RMassBank</a>	BioC
Querying information from a large number of databases.	<a href="#">webchem</a>	CRAN
R Interface to the ClassyFire REST API.	<a href="#">classyfireR</a>	CRAN
Allows mapping of identifiers from one database to another, for metabolites, genes, proteins, and interactions.	<a href="#">BridgeDbR</a>	BioC
Define utilities for exploration of human metabolome database, including functions to retrieve specific metabolite entries and data snapshots with pairwise associations.	<a href="#">hmdbQuery</a>	BioC
Parsers for many compound databases including HMDB, MetaCyc, ChEBI, FooDB, Wikidata, WikiPathways, RIKEN respect, MaConDa, T3DB, KEGG, Drugbank, LipidMaps, MetaboLights, Phenol-Explorer, MassBank.	<a href="#">MetaDBparse</a>	GitHub
Functionality to create and use compound databases generated from (mostly publicly) available resources such as HMDB, ChEBI and PubChem.	<a href="#">CompoundDb</a>	GitHub
Standardized and extensible framework to query chemical and biological databases.	<a href="#">biodb</a>	GitHub

## 2.7 Network analysis and biochemical pathways

### 2.7.1 Network infrastructure and analysis

### 2.7.2 Metabolite annotation

### 2.7.3 Generation of metabolic networks

### 2.7.4 Pathway analysis

### 2.7.5 Pathway resources and interfaces

Table 8: R packages for network analysis and Biochemical pathways.

Functionalities	Package	Repo
<b>Network infrastructure and analysis</b>		
Infrastructure for representation of networks, analysis and visualization.	<a href="#">igraph</a>	CRAN
Infrastructure for representation of networks, analysis and visualization.	<a href="#">tidygraph</a>	CRAN
Infrastructure for representation of networks, analysis and visualization.	<a href="#">statnet</a>	CRAN
Interactive visualization and manipulation of networks.	<a href="#">RedeR</a>	BioC
Comparison of correlation networks from two experiments.	<a href="#">DiffCorr</a>	CRAN
Correlation-based networks from metabolomics data and analysis tools.	<a href="#">BioNetStat</a>	BioC
<b>Annotation</b>		
Putative annotation of unknowns in MS1 data.	<a href="#">MetNet</a>	BioC
Putative annotation of unknowns in MS1 data.	<a href="#">xMSAnnotator</a>	SF
Putative annotation of unknowns using MS1 and MS2 data.	<a href="#">MetDNA</a>	GitHub
Visualization of spectral similarity networks, putative annotation of unknowns using MS2 data.	<a href="#">MetCirc</a>	BioC
Putative annotation of unknowns using MS2 data, clustering of MS2 data.	<a href="#">ChuMSID</a>	BioC
Putative annotation of unknowns using MS2 data.	<a href="#">compMS2Miner</a>	GitHub
<b>Generation of metabolite networks</b>		
Biochemical reaction networks, spectral and structural similarity networks.	<a href="#">MetaMapR</a>	GitHub
Correlation-based networks, structural similarity networks.	<a href="#">Metabox</a>	GitHub
Targeted metabolome-wide association studies.	<a href="#">MetabNet</a>	SF
Generation of scale-free correlation-based networks.	<a href="#">WGCNA</a>	CRAN
<b>Pathway analysis</b>		
Analysis of -omics data, pathway, transcription factor and target gene identification.	<a href="#">pwOmics</a>	BioC
MSEA a metabolite set enrichment analysis with factor loading in principal component analysis.	<a href="#">mseapca</a>	CRAN
Enrichment analysis of a list of affected metabolites.	<a href="#">tmod</a>	CRAN
Network-based enrichment analysis of a list of affected metabolites.	<a href="#">FELLA</a>	BioC
Pathway-based enrichment analysis of a list of affected metabolites.	<a href="#">CePa</a>	CRAN
Differential analysis, modules/sub-pathway identification using networks.	<a href="#">MetaboDiff</a>	GitHub

Table 8: R packages for network analysis and Biochemical pathways. *(continued)*

Functionalities	Package	Repo
Integrates metabolic networks and RNA-seq data to construct condition-specific series of metabolic sub-networks and applies to gene set enrichment analysis	<a href="#">metaboGSE</a>	CRAN
Differential analysis.	<a href="#">SDAMS</a>	BioC
Biomarker identification.	<a href="#">liliko</a>	CRAN
Biomarker identification.	<a href="#">INDEED</a>	BioC
Biomarker identification.	<a href="#">MoDentify</a>	GitHub
Pathway activity profiling.	<a href="#">PAPi</a>	BioC
Pathway activity profiling.	<a href="#">pathwayPCA</a>	BioC
Flux balance analysis.	<a href="#">BiGGR</a>	BioC
Flux balance analysis.	<a href="#">abcdeFBA</a>	CRAN
Flux balance analysis.	<a href="#">sybil</a>	CRAN
Flux balance analysis.	<a href="#">fbar</a>	CRAN
Identification of affected pathway from phenotype data (interface with graphite).	<a href="#">SPIA</a>	BioC
Identification of affected pathway from phenotype data (interface with graphite).	<a href="#">clipper</a>	BioC
Interface to PathVisio and WikiPathways and pathway analysis and enrichment.	<a href="#">RPathVisio</a>	GitHub
Enrichment analysis of a list of genes and metabolites.	<a href="#">RaMP</a>	GitHub
Simulation of longitudinal metabolomics data based on an underlying biological network	<a href="#">MetaboLouise</a>	CRAN
<b>Pathway resources and interfaces</b>		
BioPax parser and representation in R.	<a href="#">rBiopaxParser</a>	BioC
Interface to KEGG, Biocarta, Reactome, NCI/Nature Pathway Interaction Database, HumanCyc, Panther, SMPDB and PharmGKB.	<a href="#">graphite</a>	BioC
Interface to NCI Pathways Database.	<a href="#">NCIgraph</a>	BioC
Interface to KEGG.	<a href="#">pathview</a>	BioC
Interface to KEGG.	<a href="#">KEGGgraph</a>	BioC
Interface to systems biology markup language (SBML).	<a href="#">SBMLR</a>	BioC
Interface to systems biology markup language (SBML).	<a href="#">rsbml</a>	BioC
Interface to Gaggle-enabled software (Cytoscape, Firegoose, Gaggle Genome browser).	<a href="#">gaggle</a>	BioC
Interface to molecular interaction databases.	<a href="#">PSICQUIC</a>	BioC
Interface to KEGG REST server.	<a href="#">KEGGREST</a>	BioC
Interface to BioPAX OWL files and the Pathway Commons (PW) molecular interaction database.	<a href="#">paxtoolsr</a>	BioC
Interface to WikiPathways.	<a href="#">rWikiPathways</a>	BioC
Database that integrates metabolite and gene biological pathways from HMDB, KEGG, Reactome, and WikiPathways. Includes user-friendly R Shiny web application for queries and pathway enrichment analysis.	<a href="#">RaMP-DB</a>	GitHub

## 2.8 Multifunctional workflows

Table 9: R packages with multifunctional workflows.

Functionalities	Package	Repo
<b>NA</b>		
Convenience wrapper for pre-processing tools (XCMS, CAMERA) and a number of statistical analyses.	<a href="#">MAIT</a>	BioC
Preprocessing (XCMS), replicate merging, noise, blank and missingness filtering, feature grouping, annotation of known compounds, isotopic labeling analysis, annotation from KEGG or HMDB, common biotransformations, probabilistic putative metabolite annotation.	<a href="#">mzMatch</a>	SF
XCMS and CAMERA based workflow for non-targeted processing of LC-MS datasets, It includes pre-processing, peak picking, peak filtering, data normalization and descriptive statistics calculation.	<a href="#">MStractor</a>	GitHub
Performs simultaneous raw data to mzXML conversion (MSConvert), peak-picking, automatic PCA outlier detection and statistical analysis, visualization and possible MS2 target list determination during an MS1 metabolomic profiling experiment.	<a href="#">simExTargId</a>	GitHub
Pre-processing of large LC-MS datasets. Performs automatic PCA with iterative automatic outlier removal and, clustering analysis and biomarker discovery.	<a href="#">MetMSLine</a>	GitHub
Workflow for the systematic analysis of 1H NMR metabolomics dataset in quantitative genetics. Performs pre-processing, mQTL mapping, metabolites structural assignment and offers data visualisation tools.	<a href="#">mQTL.NMR</a>	BioC
Workflow for pre-processing, qc, annotation and statistical data analysis of LC-MS and GC-MS based metabolomics data to be submitted to public repositories.	<a href="#">MetaDB</a>	GitHub
Specmine is a framework mainly built on a number of already published packages. It supports data processing from different analytical platforms (LC-MS, GC-MS, NMR, IR, UV-Vis).	<a href="#">specmine</a>	GitHub
Common interface for a number of different MS based data processing software. It covers various aspects, such as data preparation and data extraction, formula calculation, compound identification and reporting.	<a href="#">patRoan</a>	GitHub
Processing of high resolution of LC-MS data for environmental trend analysis.	<a href="#">enviMass</a>	Zenodo
Workflow for preprocessing of LC-HRMS data, suspect screening, screening for transformation products using combinatorial prediction, and interactive filtering based on ratios between sample groups.	<a href="#">RMassScreening</a>	GitHub
Workflow to perform pre-processing, statistical analysis and metabolite identifications based on database search of detected spectra.	<a href="#">MetaboNexus</a>	GitHub
Shiny-based platform to extract differential features from LC-MS data, includes XCMS-based feature detection, statistical analysis, prediction of molecular formulas, annotation of MS2 spectra, MS2 molecular networking and chemical compound database search.	<a href="#">METABOseek</a>	GitHub
RShiny interface to Metabolomics packages & MetaboAnalyst scripts.	<a href="#">MetaboShiny</a>	GitHub

Table 9: R packages with multifunctional workflows. (*continued*)

Functionalities	Package	Repo
Preprocessing and visualizing for LC-MS data, as well as statistical analyses, mainly based on univariate linear models.	<a href="#">amp</a>	GitHub



## 2.9 User interfaces and workflow management systems

Table 10: Packages to interface R with other languages and workflow environments

Functionalities	Package	Repo
<b>NA</b>		
Given an R function and its manual page, make the documented function available in Galaxy.	<a href="#">RGalaxy</a>	BioC
Integration of R and C++. Many R data types and objects can be mapped back and forth to C++ equivalents.	<a href="#">Rcpp</a>	CRAN
Low-Level R to Java Interface.	<a href="#">rJava</a>	CRAN
Interface to 'Python' modules, classes, and functions and translation between R and Python objects.	<a href="#">reticulate</a>	CRAN

## 2.10 Metabolomics data sets

Table 11: Metabolomics data sets packaged as R packages.

Functionalities	Package	Repo
<b>LC-MS</b>		
12 HPLC-MS NetCDF files (Agilent 1100 LC-MSD SL).	<a href="#">faahKO</a>	BioC
16 UPLC-MS mzData files (Bruker microTOFq).	<a href="#">mtbls2</a>	BioC
12 UPLC-MS mzML files (AB Sciex TripleTOF 5600, SWATH mode).	<a href="#">mtbls297</a>	GitHub
Different raw MS files (LTQ, TripleQ, FTICR, Orbitrap, QTOF) some in different formats (mzML, mzXML, mzData, mzData.gz, NetCDF, mz5). Also mzid format from proteomics.	<a href="#">msdata</a>	BioC
Metadata and DDA MS/MS spectra of 15 narcotics standards (LTQ Orbitrap XL).	<a href="#">RMassBankData</a>	BioC
183 x 109 peak table.	<a href="#">ropls</a>	BioC
69 x 5,501 peak table.	<a href="#">biosigner</a>	BioC
40 x 1,632 peak table.	<a href="#">BioMark</a>	CRAN
Raw MS files from a set of blanks and standards that contain common environmental contaminants (acquired with Bruker maXis 4G).	<a href="#">patRoonaData</a>	GitHub
Proteomics, metabolomics GC-MS and Lipidomics data from Calu-3 cell culture; 3 mockulum treated and 9 MERS-CoV treated; Time point, 18 hour from MassIVE dataset ids MSV000079152, MSV000079153, MSV000079154.	<a href="#">pmartRdata</a>	GitHub
<b>FIA-MS</b>		
6 mzML files (human plasma spiked with 40 compounds acquired in positive mode on an orbitrap fusion).	<a href="#">plasFIA</a>	BioC
mzML files (Thermo Exactive) from comparison of leaf tissue from 4 B. distachyon ecotypes with Flow-infusion electrospray ionisation-high resolution mass spectrometry (FIE-HRMS). Also includes data sets with 10 technical injections of human urine and another 10 injections from leaf tissue (ecotype ABR1).	<a href="#">metaboData</a>	GitHub
<b>GC-MS</b>		
52 x 154 peak table.	<a href="#">pcaMethods</a>	BioC
<b>NMR</b>		
18 x 189 peak table.	<a href="#">MetabolAnalyze</a>	CRAN
33 x 164 peak table.	<a href="#">MetabolAnalyze</a>	CRAN
ASICSdata: 1D NMR spectra for ASICS.	<a href="#">ASICSdata</a>	BioC

### 3 Conclusions

### 4 Appendices

## Appendix 1: The MSP File Format and package support

---

```
Name: unknown
Num Peaks: 2
85.345 100; 76.321 50;
```

---

**Listing S1:** Example for the basic NIST format.

Name: 1-Methylhistidine Synon: (2S)-2-amino-3-(1-methyl-1H-imidazol-4-yl)propanoic acid SYNON: \$.00in-source DB#: HMDB0000001_c_ms_1469 InChIKey: BRMWTNUJHUMWMS-LURJTMIESA-N Instrument_type: GC-MS Retention_index: 1807.71 Formula: C7H11N3O2 MW: 169 ExactMass: 169.0851 Comments: "column=5%-phenyl-95%-dimethylpolysiloxane capillary column" "derivatization type=2 TMS" "derivatization formula=C13H27N3O2Si2" "derivative mw=313.544" "retention index=1807.71" "retention index type=based on 9 n-alkanes (C10-C36)" "instrument type=GC-MS" "chromatography type=GC" "cas number=332-80-9" "molecular formula=C7H11N3O2" "total exact mass=169.085126592" "InChIKey=BRMWTNUJHUMWMS-LURJTMIESA-N" Num Peaks: 10 70 0.014; 71 0.007; 72 0.02; 76 0.008; 77 0.008; 78 0.002; 79 0.003; 80 0.005; 81 0.108; 82 0.017;	NAME: Aspartame; LC-ESI-ITFT; MS2; CE PRECURSORMZ: 295.128848 PRECURSORTYPE: [M+H] <sup>+</sup> INSTRUMENTTYPE: LC-ESI-ITFT SMILES: COC(=O)C(CC1=CC=CC=C1)N=C(O)C(N)CC(O)=O INCHIKEY: IAOZJIPTCAWIRG-UHFFFAOYNA-N Ontology: Peptides COLLISIONENERGY: 35 FORMULA: C14H18N2O5 RETENTIONTIME: IONMODE: Positive Comment: registered in MassBank Num Peaks: 9 120.0804 13 180.10201 138 217.0968 14 235.10789 390 245.0921 274 260.09171 132 263.1026 286 277.11859 1000 278.1022 28
<b>Listing S2:</b> Example for the canonical NIST format.	<b>Listing S3:</b> RIKEN PRIME msp format example.

**Table S1:** Overview of MS/MS handling in different R packages. ‘-’ means not available, for the remaining entries see the text above.

package	read msp	write msp	spectral matching and additional information
baitmet			N vs DB; cosine, Stein & Scott composite similarity product
compMS2Miner	NIST, RIKEN PRIME msp	RIKEN PRIME msp	N vs DB; dot product
enviGCMS		basic NIST	
erah	NIST	only result export	N vs DB; cosine
flagme		only result export	
metaMS	NIST	NIST; slow	1 vs DB, N vs DB; proprietary
MatchWeiz			N vs DB; X-Rank
MetCirc			N vs N; normalized dot product; will switch to MSnbase functions soon
MSeasy		only result export	N vs DB; Queries the NIST mass spectral search tool
MSnbase	**	**	1 vs 1, N vs N; dot product and more, user def.
msPurity			N vs DB; dot product
OrgMassSpecR	basic NIST	basic NIST	1 vs 1; normalized dot product
RAMClustR			RAMClustR can import and utilize spectrum similarities from MS-FINDER;
rTANDEM			N vs DB; dot product; R wrapper for X!Tandem software
SwathXtend-	(PeakView / OpenSWATH)	-(PeakView / OpenSWATH)	
TargetSearch	NIST (with error)	NIST	N vs DB; RI-based

## Appendix 2: metaRbolomics dependencies network

### Libraries and settings

```
options("repos" = list(CRAN="http://cran.rstudio.com/"))

library(devtools)    # for revdep()
library(igraph)      # for graph_from_edgelist/( and simplify() )
library(visNetwork)  # for visNetwork() and friends
library(networkD3)   # for saveNetwork()
library(chromote)     # for default_chromote_object()
library(webshot2)    # for webshot()
library(png)         # For displaying an image
library(dplyr)
library(purrr)

source("revDepNetHelper.R")

set_default_chromote_object(Chromote$new(browser = Chrome$new(args = "--no-sandbox")))
```

### Read package names from our table

```
reviewTables <- read.delim("data/AllMetaRbolomicsTables.csv", stringsAsFactors = FALSE)
reviewPkgs <- reviewTables[, "Package"]

pkgs <- reviewPkgs
```

### Get reverse dependencies

#### 4.0.0.1 For CRAN and BioC packages

```
e1 <- sapply(pkgs, function(pkg) {
  rd <- revdep(pkg, dependencies = c("Depends", "Imports", "LinkingTo"),
    recursive = FALSE, ignore = NULL, bioconductor = TRUE)
  as.matrix(cbind(Package=rep(pkg, length.out=length(rd)), ReverseDep=rd))
})
e1 <- do.call(rbind, e1)
```

#### 4.0.0.2 For GitHub and GitLab

The above `devtools::revdep` cannot read from GitHub/GitLab repositories. We have a helper function that downloads and parses the DESCRIPTION file from GitHub/GitLab. Since we cannot get reverse dependencies directly for GitHub/GitLab packages, those packages they are only used as additional reverse dependencies for the CRAN/BioC packages.

```
gitdeps_reverse <- reviewTables %>%
  mutate(dep_tree = map(Code_link, get_git_deps)) %>%
  pull(dep_tree) %>%
  bind_rows() %>%
```

```

filter(Dep %in% el[, "Package"]) %>%
rename(Package = Dep, ReverseDep = Package) %>%
as.matrix()

```

```

## Warning in readLines(file): incomplete final line found on 'C:
## \Users\jan\AppData\Local\Temp\RtmpEDPsdD\file260840ca1fcb'

```

```

## Warning in readLines(file): incomplete final line found on 'C:
## \Users\jan\AppData\Local\Temp\RtmpEDPsdD\file260862906c1a'

```

```

el <- rbind(el, gitdeps_reverse)

```

## Building dependency network

In total, we were analysing 288 packages. For each package, this returns the set of packages in CRAN or BioC that depend on, import from or link to the package (i.e., its direct reverse dependencies) using the `devtools::revdep()` function. A few packages with the highest number of reverse dependencies have been excluded, as they would dominate the visualisation. It was not possible to detect reverse dependencies from other hosting places such as GitHub or GitLab.

From the total, 65 packages had at least one such reverse dependency.

```

## Remove packages with most reverse dependencies
## which would dominate the network

el <- el[! el[, "Package"] %in% c("Rcpp", "igraph", "vegan", "caret", "rJava", "reticulate"), ]

## Create graph, and simplify redundancy
g <- graph_from_edgelist(el, directed = TRUE)
g <- igraph::simplify(g, remove.multiple = TRUE, remove.loops = TRUE)

# get data and plot :
data <- toVisNetworkData(g)

data$nodes <- cbind(data$nodes,
                    font.size=30,
                    color.background = ifelse(data$nodes[, "id"] %in% pkgs ,
                                              rgb(0, 0, 200, 128, max = 255),
                                              rgb(0, 200, 0, 128, max = 255)))

vn <- visNetwork(nodes = data$nodes,
                 edges = data$edges,
                 width=1000, height=1000) %>%
  visPhysics(timestep = 0.3,
            barnesHut = list(centralGravity=0.35,
                             springLength = 95)) %>%
  visOptions(highlightNearest = TRUE)

vn

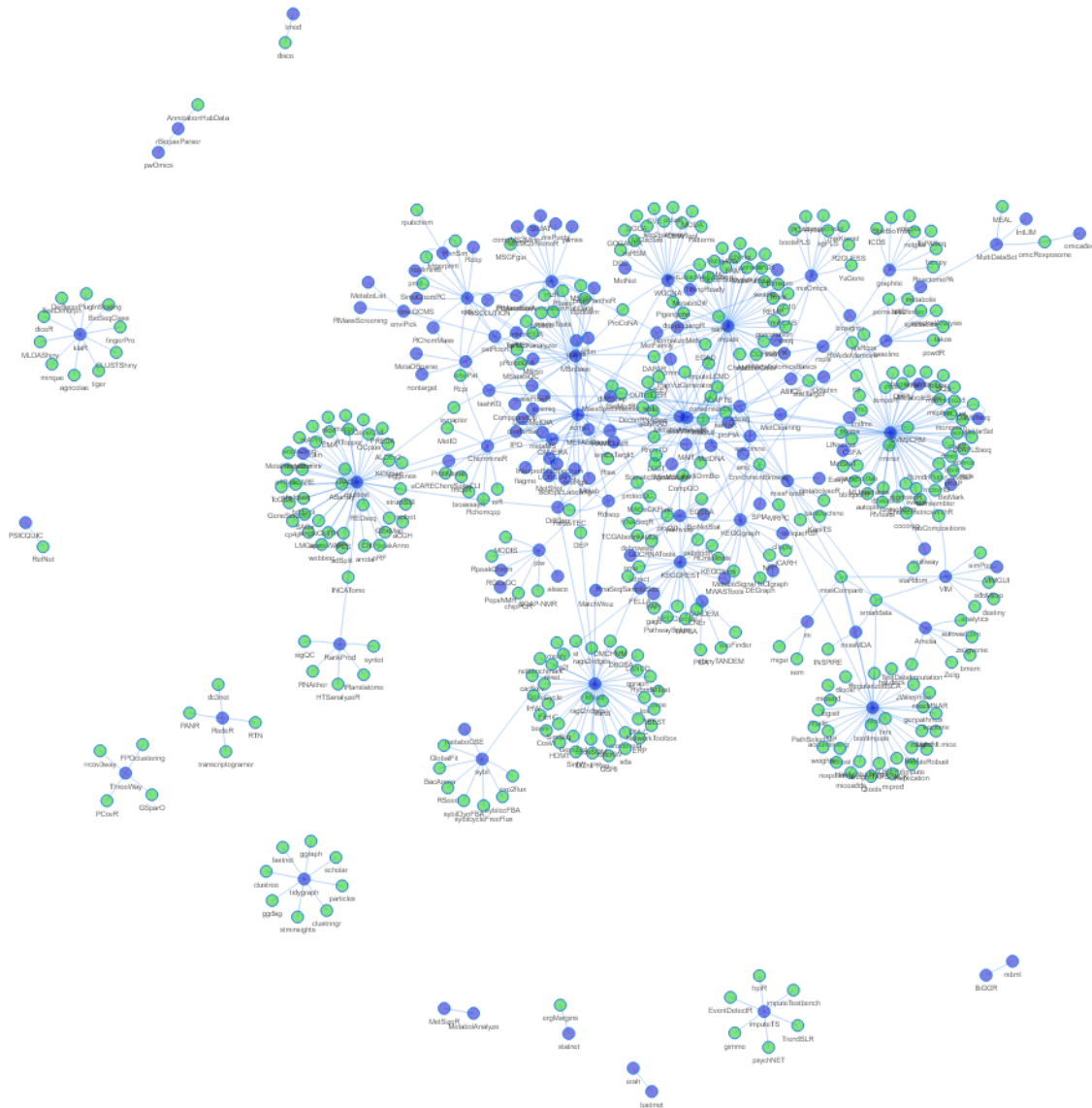
```

Figure S1: Dependency network of R packages. Shown in blue are packages mentioned in the review. Edges connect to packages that depend on another package, as long as that is in CRAN or BioC. Green nodes correspond to packages in CRAN or BioC not covered in the review. Not shown are 1) infrastructure packages e.g. rJava, Rcpp 2) packages from the review without reverse dependencies and 3) data packages. Some packages from the review are not in current versions of CRAN or BioC. An interactive version of this figure is available from <https://stanstrup.gitlab.io/metaRbolomics-book/appendix-2-metarbolomics-dependencies-network.html>.



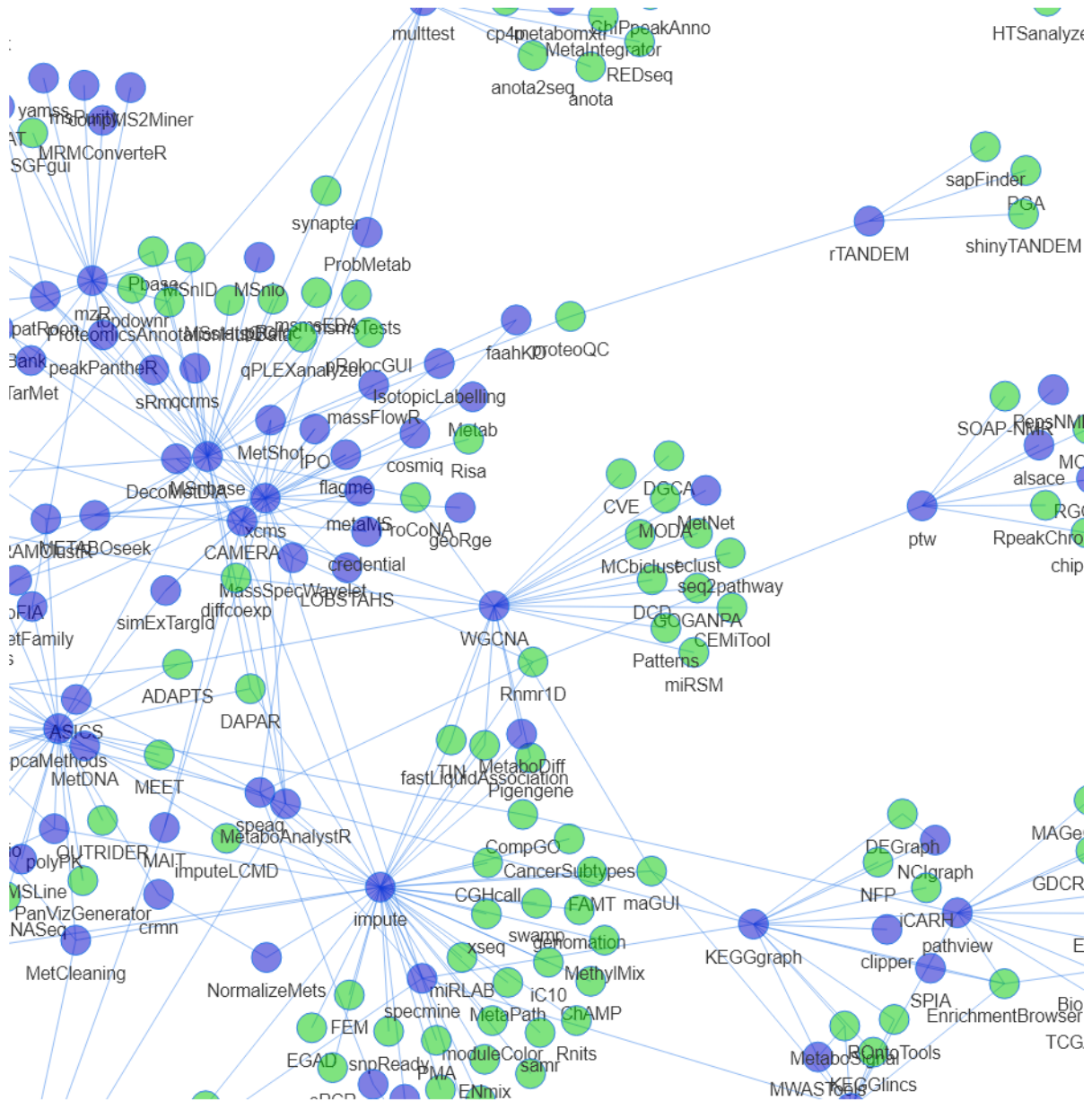
## Save network plot

```
saveNetwork(vn, "vn.html")
webshot("vn.html", "revDepNet-60.png", delay = 60)
```



```
vnZoom <- visNetwork(nodes = data$nodes,
  edges = data$edges,
  width=1000, height=1000) %>%
```

```
visIgraphLayout()%>%
visEvents(type="once", startStabilizing = 'function() {
  this.fit({nodes:["ptw", "Rnmr1D", "RpeakChrom", "alsace",
    "PepsNMR", "ASICS", "MODIS", "RGCxGC"]})
})
saveNetwork(vnZoom, "vnZoom.html")
webshot("vnZoom.html", "revDepNet-zoom.png", delay = 5)
```



You can access the files at:

- [vn.html](#)
- [revDepNet-60.png](#)
- [vnZoom.html](#)
- [revDepNet-zoom.png](#)

## Notes

The source code for this page is on GitHub at [gitlab.com/stanstrup/metaRbolomics-book](https://gitlab.com/stanstrup/metaRbolomics-book)

The HTML output is shown at <https://stanstrup.gitlab.io/metaRbolomics-book/appendix-2-metarbolomics-dependencies-network.html>

and <https://stanstrup.gitlab.io/metaRbolomics-book/vn.html> (Caveat: long rendering time, blank page without any visible progress)

This page was created with the following packages:

```
sessionInfo()
```

```
## R version 3.5.2 (2018-12-20)
## Platform: x86_64-w64-mingw32/x64 (64-bit)
## Running under: Windows >= 8 x64 (build 9200)
##
## Matrix products: default
##
## locale:
## [1] LC_COLLATE=English_United States.1252
## [2] LC_CTYPE=English_United States.1252
## [3] LC_MONETARY=English_United States.1252
## [4] LC_NUMERIC=C
## [5] LC_TIME=English_United States.1252
##
## attached base packages:
## [1] stats      graphics  grDevices  utils      datasets  methods   base
##
## other attached packages:
## [1] magrittr_1.5          miniCRAN_0.2.12      desc_1.2.0
## [4] png_0.1-7            webshot2_0.0.0.9000  chromote_0.0.0.9001
## [7] networkD3_0.4        visNetwork_2.0.8     igraph_1.2.4.1
## [10] devtools_2.1.0       usethis_1.5.1        tikzDevice_0.12.3
## [13] purrr_0.3.2          kableExtra_1.1.0.9000 dplyr_0.8.3
## [16] DT_0.8               readr_1.3.1          knitr_1.24
##
## loaded via a namespace (and not attached):
## [1] httr_1.4.1           pkgload_1.0.2        jsonlite_1.6
## [4] viridisLite_0.3.0    shiny_1.3.2          assertthat_0.2.1
## [7] BiocManager_1.30.4   yaml_2.2.0           remotes_2.1.0
## [10] sessioninfo_1.1.1    pillar_1.4.2         backports_1.1.4
## [13] glue_1.3.1           digest_0.6.20        promises_1.0.1.9002
## [16] rvest_0.3.4          colorspace_1.4-1     websocket_1.1.0
## [19] htmltools_0.3.6      httpuv_1.5.1         XML_3.98-1.20
## [22] pkgconfig_2.0.2      bookdown_0.13        xtable_1.8-4
## [25] scales_1.0.0         webshot_0.5.1        processx_3.4.1
## [28] later_0.8.0.9003     tibble_2.1.3         withr_2.1.2
```

## [31] cli_1.1.0	crayon_1.3.4	mime_0.7
## [34] memoise_1.1.0	evaluate_0.14	ps_1.3.0
## [37] fs_1.3.1	xml2_1.2.2	pkgbuild_1.0.5
## [40] tools_3.5.2	prettyunits_1.0.2	hms_0.5.1
## [43] stringr_1.4.0	munsell_0.5.0	callr_3.3.1
## [46] compiler_3.5.2	rlang_0.4.0	grid_3.5.2
## [49] rstudioapi_0.10	htmlwidgets_1.3	filehash_2.4-2
## [52] crosstalk_1.0.0	rmarkdown_1.15	testthat_2.2.1
## [55] codetools_0.2-15	curl_4.0	R6_2.4.0
## [58] fastmap_1.0.0	zeallot_0.1.0	rprojroot_1.3-2
## [61] stringi_1.4.3	Rcpp_1.0.2	vctrs_0.2.0
## [64] tidyselect_0.2.5	xfun_0.9	