

Package ‘MsDataHub’

April 1, 2025

Title Mass Spectrometry Data on ExperimentHub

Version 1.6.0

Description The MsDataHub package uses the ExperimentHub infrastructure to distribute raw mass spectrometry data files, peptide spectrum matches or quantitative data from proteomics and metabolomics experiments.

License Artistic-2.0

BugReports <https://github.com/RforMassSpectrometry/MsDataHub/issues>

URL <https://rformassspectrometry.github.io/MsDataHub>

Imports ExperimentHub, utils

Suggests ExperimentHubData, DT, BiocStyle, knitr, rmarkdown, testthat (>= 3.0.0), Spectra, mzR, PSMatch, QFeatures (>= 1.13.3)

biocViews ExperimentHubSoftware, MassSpectrometry, Proteomics, Metabolomics

Encoding UTF-8

VignetteBuilder knitr

Roxygen list(markdown = TRUE)

RoxygenNote 7.3.1

Config/testthat.edition 3

git_url <https://git.bioconductor.org/packages/MsDataHub>

git_branch RELEASE_3_20

git_last_commit bceb0f3

git_last_commit_date 2024-10-29

Repository Bioconductor 3.20

Date/Publication 2025-03-31

Author Laurent Gatto [aut, cre] (<<https://orcid.org/0000-0002-1520-2268>>), Kristina Gomoryova [ctb] (<<https://orcid.org/0000-0003-4407-3917>>), Johannes Rainer [aut] (<<https://orcid.org/0000-0002-6977-7147>>)

Maintainer Laurent Gatto <laurent.gatto@uclouvain.be>

Contents

benchmarkingDIA	2
cdf	2
cptac	3
MsDataHub	4
PXD000001	4
Report.Derks2022.plexDIA	5
sciex	5
TripleTOF	6
Index	7

benchmarkingDIA	<i>DIA benchmarking data</i>
-----------------	------------------------------

Description

These data were generated based on publicly available DIA benchmarking dataset from Gotti et al. (2021). A subset of raw data, containing "overlapped" in the File.Name were searched using the DIA-NN software, and the resulting report.tsv (here labelled as 'benchmarkingDIA.tsv') is provided.

The dataset contains 8 conditions containing a mix of E.coli and Universal Standard Protein-1 (UPS1) peptides. Per 1 ug of E.coli protein (equal in all samples), UPS1 proteins are diluted to final concentration of 50, 25, 10, 5, 2.5, 1, 0.25 and 0.1 fmol.

Each sample was prepared in 3 replicates, so altogether there are 24 samples in the dataset.

Author(s)

Kristina Gomoryova and Laurent Gatto

References

- Gotti C, Roux-Dalvai F, Joly-Beauparlant C, Mangnier L, Leclercq M, Droit A. Extensive and Accurate Benchmarking of DIA Acquisition Methods and Software Tools Using a Complex Proteomic Standard. *J Proteome Res.* 2021 Oct 1;20(10):4801-4814. doi: 10.1021/acs.jproteome.1c00490. Epub 2021 Sep 2. PMID: 34472865.

cdf	<i>MS data in CDF format</i>
-----	------------------------------

Description

This data set represents a single CDF file in (AIA/ANDI) NetCDF format from a larger experiment in which the metabolic consequences of knocking the fatty acid amide hydrolase (FAAH) gene in mice was investigated. The file contains data in centroid mode acquired in positive ion mode from 200-600 m/z and 2500-4500 seconds.

Data file:

- ko15.CDF* file in NetCDF format.

References

- Saghatelian, A et al. *Assignment of endogenous substrates to enzymes by global metabolite profiling*, Biochemistry, 2004. <http://dx.doi.org/10.1021/bi0480335>

cptac

CPTAC label-free data

Description

This case-study is a subset of the data of the 6th study of the Clinical Proteomic Technology Assessment for Cancer (CPTAC) (Paulovich et al. 2010). In this experiment, the authors spiked the Sigma Universal Protein Standard mixture 1 (UPS1) containing 48 different human proteins in a protein background of 60 micro g/micro L *Saccharomyces cerevisiae* strain BY4741.

Five different spike-in concentrations were used:

- 6A: 0.25 fmol UPS1 proteins/micro L
- 6B: 0.74 fmol UPS1 proteins/micro L
- 6C: 2.22 fmol UPS1 proteins/micro L
- 6D: 6.67 fmol UPS1 proteins/micro L
- 6E: 20 fmol UPS1 proteins/micro L

Three replicates are available for each concentration.

The data were searched with MaxQuant version 1.5.2.8 (Cox et al. 2008) including matching between runs. Detailed search settings were described in Goeminne et al. (2016).

Three files are readily available as tab-delimited spreadsheets:

- cptac_a_b_peptides.txt: triplicates from lab 3 for groups 6A and 6B.
- cptac_a_b_c_peptides.txt: triplicates from labs 1, 2 and 3 for groups 6A, 6B and 6C.
- cptac_peptides.txt: triplicates from labs 1, 2, and 3 for all groups.

Author(s)

Laurent Gatto and Lieven Clement

References

- Paulovich, Amanda G, Dean Billheimer, Amy-Joan L Ham, Lorenzo Vega-Montoto, Paul A Rudnick, David L Tabb, Pei Wang, et al. 2010. *Interlaboratory Study Characterizing a Yeast Performance Standard for Benchmarking LC-MS Platform Performance*. Mol. Cell. Proteomics 9 (2): 242–54.
- Cox, J, and M Mann. 2008. *MaxQuant Enables High Peptide Identification Rates, Individualized p.p.b.-Range Mass Accuracies and Proteome-Wide Protein Quantification*. Nat Biotechnol 26 (12): 1367–72. <https://doi.org/10.1038/nbt.1511>.
- Goeminne, LJ, Gevaert K and Clement, L. 2016. *Peptide-level Robust Ridge Regression Improves Estimation, Sensitivity, and Specificity in Data-dependent Quantitative Label-free Shotgun Proteomics*, Mol Cell Proteomics, 15:2 657-668.

MsDataHub	All MsDataHub datasets
-----------	------------------------

Description

The MsDataHub package provides example mass spectrometry data, peptide spectrum matches or quantitative data from proteomics and metabolomics experiments.

The `MsDataHub()` function returns a `data.frame` with all the annotated datasets provided in the package. For details on these individual datasets, refer to their respective manual pages.

See the vignette and the respective manuals pages for more details about the package and the data themselves.

Usage

```
MsDataHub()
```

Value

A `data.frame` describing the data available in MsDataHub.

Author(s)

Laurent Gatto

Examples

```
MsDataHub()
```

PXD000001	PXD000001 Proteomics Data
-----------	---------------------------

Description

The PXD000001 files are part of the first ProteomeXchange submission (Vizcaíno J.A. et al, 2014), and contain the following files.

- TMT_Erwinia_1uLSike_Top10HCD_isol2_45stepped_60min_01-20141210.mzML.gz: an TMT6-plex LC-MSMS data containing 6 human spiked-in proteins in a constant *Erwinia carotovora* protein background. The data is described in more details in Gatto and Christoforou (2013).
- TMT_Erwinia_1uLSike_Top10HCD_isol2_45stepped_60min_01-20141210.mzid: generated searching the raw data against the *Erwinia carotovora* fasta database

References

- Vizcaíno J.A. et al. *ProteomeXchange: globally co-ordinated proteomics data submission and dissemination*, Nature Biotechnology 2014, 32, 223–226. <http://www.ncbi.nlm.nih.gov/pubmed/24727771>
- Gatto L. and Christoforou A. *Using R and Bioconductor for proteomics data analysis*, Biochim Biophys Acta - Proteins and Proteomics, 2013. <http://www.ncbi.nlm.nih.gov/pubmed/23692960>

See Also

The `rpx` package can be used to access and download any PRIDE/ProteomeXchange files.

Report.Derks2022.plexDIA
Derks 2022 plexDIA data

Description

Single cell proteomics data acquired by the Slavov Lab using the plexDIA protocol. It contains quantitative information from pancreatic ductal acinar cells (PDAC; HPAF-II), melanoma cells (WM989-A6-G3) and monocytes (U-937) at precursor and protein level. The each run acquired 3 samples thanks to mTRAQ multiplexing.

The data were downloaded from the Slavov lab google drive:

- https://drive.google.com/drive/folders/1pUC2zgXKtKYn22mlor0lmUDK0frgwL_-
- DIANN_outputs
- wJD1146_1193_1200_tsvLib
- Report.tsv

For more details about the data: <https://plexdia.slavovlab.net/>

The file is reshare here allow its dissemination via the MsDataHub package.

Author(s)

Laurent Gatto

References

Derks, J., Leduc, A., Wallmann, G. et al. Increasing the throughput of sensitive proteomics by plexDIA. Nat Biotechnol (2022). 10.1038/s41587-022-01389-w.

sciex *AB Sciex LC-MS data files*

Description

The `sciex` mzML files represent profile-mode LC-MS data of pooled human serum samples (the same pool being measured). The samples were analyzed by ultra high-performance liquid chromatography (UHPLC; Agilent 1290) coupled to a Q-TOF mass spectrometer (TripleTOF 5600+ AB Sciex). The chromatographic separation was based in hydrophilic interaction liquid chromatography (HILIC) and performed using an Waters Acquity BEH Amide, 100 x 2.1 mm column.

The mass spectrometer was operated in full scan mode in the mass range from 50 to 1000 m/z and with an accumulation time of 250 ms. The files represent a subset of spectra/scans from m/z 105 to 134 and from retention time 0 to 260 seconds. The files were generated in the same LC-MS run, but from different injections. Details on the individual files are provided below.

Files:

- *20171016_POOL_POS_1_105-134.mzML*: profile-mode LC-MS data of pooled human serum samples. Injection index: 1.
- *20171016_POOL_POS_3_105-134.mzML*: profile-mode LC-MS data of pooled human serum samples. Injection index: 19.

Author(s)

Sigurdur Smarason, Giuseppe Paglia and Johannes Rainer

TripleTOF

Triple TOF SWATH Data

Description

These files represent data from reverse-phased LC-MS/MS runs on the Agilent Pesticide mix obtained from a Sciex 6600 Triple ToF operated either in Sequential Window Acquisition of all Theoretical mass spectra (SWATH) or Data Dependent Acquisition (DDA) acquisition mode.

The data files are:

- *PestMix1_DDA.mzML*: mzML file with MS1 and MS2 spectra from the Agilent Pesticide Mix acquired in DDA mode.
- *PestMix1_SWATH.mzML*: mzML file with MS1 and MS2 spectra from the Agilent Pesticide Mix acquired in SWATH mode.

Author(s)

Micheal Witting, Johannes Rainer

Index

20171016_POOL_POS_1_105-134.mzML
(sciex), 5
20171016_POOL_POS_3_105-134.mzML
(sciex), 5

benchmarkingDIA, 2

cdf, 2
cptac, 3
cptac_a_b_c_peptides.txt (cptac), 3
cptac_a_b_peptides.txt (cptac), 3
cptac_peptides.txt (cptac), 3

ko15.CDF (cdf), 2

MsDataHub, 4
MsDataHub(), 4

PestMix1_DDA.mzML (TripleTOF), 6
PestMix1_SWATH.mzML (TripleTOF), 6
PXD000001, 4

Report.Derks2022.plexDIA, 5

sciex, 5

TMT_Erwinia_1uLSike_Top10HCD_isol2_45stepped_60min_01-20141210.mzid
(PXD000001), 4
TMT_Erwinia_1uLSike_Top10HCD_isol2_45stepped_60min_01-20141210.mzML.gz
(PXD000001), 4
TMT_Erwinia_1uLSike_Top10HCD_isol2_45stepped_60min_01.20141210.mzid
(PXD000001), 4
TMT_Erwinia_1uLSike_Top10HCD_isol2_45stepped_60min_01.20141210.mzML.gz
(PXD000001), 4

TripleTOF, 6

X20171016_POOL_POS_1_105.134.mzML
(sciex), 5
X20171016_POOL_POS_3_105.134.mzML
(sciex), 5