

Package ‘RaMS’

December 14, 2022

Type Package

Title R Access to Mass-Spec Data

Version 1.3.1

Maintainer William Kumler <wkumler@uw.edu>

Description R-based access to mass-spectrometry (MS) data. While many packages exist to process MS data, many of these make it difficult to access the underlying mass-to-charge ratio (m/z), intensity, and retention time of the files themselves. This package is designed to format MS data in a tidy fashion and allows the user perform the plotting and analysis.

License MIT + file LICENSE

URL <https://github.com/wkumler/RaMS>

BugReports <https://github.com/wkumler/RaMS/issues>

Encoding UTF-8

Imports xml2, base64enc, data.table, utils

RoxygenNote 7.2.1

Suggests testthat, knitr, rmarkdown, tidyverse, ggplot2, dplyr, plotly, openxlsx, DBI, RSQLite, reticulate

VignetteBuilder knitr

NeedsCompilation no

Author William Kumler [aut, cre, cph],
Ricardo Cunha [ctb],
Ethan Bass [ctb]

Repository CRAN

Date/Publication 2022-12-14 21:40:06 UTC

R topics documented:

checkOutputQuality	2
getEncoded	3

giveEncoding	4
grabAccessionData	4
grabMSdata	5
grabMzmlBPC	7
grabMzmlDAD	8
grabMzmlData	8
grabMzmlEncodingData	10
grabMzmlMetadata	11
grabMzmlMS1	11
grabMzmlMS2	12
grabMzxmlBPC	12
grabMzxmlData	13
grabMzxmlEncodingData	15
grabMzxmlMetadata	15
grabMzxmlMS1	16
grabMzxmlMS2	16
grabMzxmlSpectraMzInt	17
grabMzxmlSpectraPremz	17
grabMzxmlSpectraRt	18
grabMzxmlSpectraVoltage	18
grabSpectraInt	19
grabSpectraMz	19
grabSpectraPremz	20
grabSpectraRt	20
grabSpectraVoltage	21
minifyMSdata	21
minifyMzml	23
minifyMzxml	24
msdata_connection	26
node2dt	26
pmppm	27
print.msdata_connection	27
tmzmlMaker	28
[.msdata_connection	29
\$.msdata_connection	30

Index 31

checkOutputQuality	<i>Check that the output data is properly formatted.</i>
--------------------	--

Description

This function checks that data produced by repeated calls to the ‘grabMzmlData()’ and ‘grabMzxmlData()’ functions is formatted properly before it’s provided to the user. It checks that all of the requested data has been obtained and warns if data is found to be empty, misnamed, or has columns of the wrong type.

Usage

```
checkOutputQuality(output_data, grab_what)
```

Arguments

output_data	The collected data resulting from repeated calls to 'grabMzmlData()', after being bound together.
grab_what	The names of the data requested by the user.

Value

NULL (invisibly). The goal of this function is its side effects, i.e. throwing errors and providing info when the files are not found.

getEncoded	<i>Convert from compressed binary to R numeric vector</i>
------------	---

Description

Convert from compressed binary to R numeric vector

Usage

```
getEncoded(mzint_nodes, compression_type, bin_precision, endi_enc)
```

Arguments

mzint_nodes	The XML nodes containing the compressed binary string
compression_type	Compression type to be used by memDecompress
bin_precision	The bit (?) precision used by readBin
endi_enc	The byte order (?) of the string. For mzML this is always "little" but mzXML can also be "big"

Value

A numeric vector of m/z or intensity values

giveEncoding	<i>Convert from R numeric vector to compressed binary</i>
--------------	---

Description

Convert from R numeric vector to compressed binary

Usage

```
giveEncoding(mzint_vals, compression_type, bin_precision, endi_enc)
```

Arguments

mzint_vals	A numeric vector of m/z or intensity values
compression_type	Compression type to be used by memCompress
bin_precision	The bit (?) precision used by writeBin
endi_enc	The byte order (?) of the string. For mzML this is always "little" but mzXML can also be "big"

Value

A single base64-encoded string of compressed binary values

grabAccessionData	<i>Get arbitrary metadata from an mzML file by accession number</i>
-------------------	---

Description

Get arbitrary metadata from an mzML file by accession number

Usage

```
grabAccessionData(filename, accession_number)
```

Arguments

filename	The name of the file for which metadata is requested. Both absolute and relative paths are acceptable.
accession_number	The HUPO-PSI accession number for the metadata to be extracted. These accession numbers are typically of the form MS:##### and the full list can be found and searched at https://raw.githubusercontent.com/HUPO-PSI/psi-ms-CV/master/psi-ms.obo .

Value

A data frame with the name and value of the parameter requested, as deduced from the XML tag attributes corresponding to the accession number.

Examples

```
library(RaMS)
sample_dir <- system.file("extdata", package = "RaMS")
sample_file <- list.files(sample_dir, full.names=TRUE)[3]
# Get ion injection time
iit_df <- grabAccessionData(sample_file, "MS:1000927")
# Manually create TIC
int_df <- grabAccessionData(sample_file, "MS:1000285")
rt_df <- grabAccessionData(sample_file, "MS:1000016")
tic <- data.frame(rt=rt_df$value, int=int_df$value)
plot(tic$rt, tic$int, type = "l")
```

grabMSdata

Grab mass-spectrometry data from file(s)

Description

The main ‘RaMS’ function. This function accepts a list of the files that will be read into R’s working memory and returns a list of ‘data.table’s containing the requested information. What information is requested is determined by the ‘grab_what’ argument, which can include MS1, MS2, BPC, TIC, or metadata information. This function serves as a wrapper around both ‘grabMzmlData’ and ‘grabMzxmlData’ and handles multiple files, but those two have also been exposed to the user in case super-simple handling is desired. Retention times are reported in minutes, and will be converted automatically if they are encoded in seconds.

Usage

```
grabMSdata(
  files,
  grab_what = "everything",
  verbosity = NULL,
  mz = NULL,
  ppm = NULL,
  rtrange = NULL,
  prefilter = -1
)
```

Arguments

files A character vector of filenames to read into R’s memory. Both absolute and relative paths are acceptable.

grab_what	What data should be read from the file? Options include "MS1" for data only from the first spectrometer, "MS2" for fragmentation data, "BPC" for rapid access to the base peak chromatogram, "TIC" for rapid access to the total ion chromatogram, "DAD" for DAD (UV) data, and "chroms" for precompiled chromatogram data (especially useful for MRM but often contains BPC/TIC in other files). Metadata can be accessed with "metadata", which provides information about the instrument and time the file was run. These options can be combined (i.e. 'grab_data=c("MS1", "MS2", "BPC")') or this argument can be set to "everything" to extract all of the above. Options "EIC" and "EIC_MS2" are useful when working with files whose total size exceeds working memory - it first extracts all relevant MS1 and MS2 data, respectively, then discards data outside of the mass range(s) calculated from the provided mz and ppm. The default, "everything", includes all MS1, MS2, BPC, TIC, and metadata.
verbosity	Three levels of processing output to the R console are available, with increasing verbosity corresponding to higher integers. A verbosity of zero means that no output will be produced, useful when wrapping within larger functions. A verbosity of 1 will produce a progress bar using base R's txtProgressBar function. A verbosity of 2 or higher will produce timing output for each individual file read in. The default, NULL, will select between 1 and 2 depending on the number of files being read: if a single file, verbosity is set to 2; if multiple files, verbosity is set to 1.
mz	A vector of the mass-to-charge ratio for compounds of interest. Only used when combined with 'grab_what = "EIC"' (see above). Multiple masses can be provided.
ppm	A single number corresponding to the mass accuracy (in parts per million) of the instrument on which the data was collected. Only used when combined with 'grab_what = "EIC"' (see above).
rtrange	Only available when parsing mzML files. A vector of length 2 containing an upper and lower bound on retention times of interest. Providing a range here can speed up load times (although not enormously, as the entire file must still be read) and reduce the final object's size.
prefilter	A single number corresponding to the minimum intensity of interest in the MS1 data. Data points with intensities below this threshold will be silently dropped, which can dramatically reduce the size of the final object. Currently only works with MS1 data, but could be expanded easily to handle more.

Value

A list of 'data.table's, each named after the arguments requested in grab_what. E.g. \$MS1 contains MS1 information, \$MS2 contains fragmentation info, etc. MS1 data has four columns: retention time (rt), mass-to-charge (mz), intensity (int), and filename. MS2 data has six: retention time (rt), precursor m/z (premsz), fragment m/z (fragmsz), fragment intensity (int), collision energy (voltage), and filename. Data requested that does not exist in the provided files (such as MS2 data requested from MS1-only files) will return an empty (length zero) data.table. The data.tables extracted from each of the individual files are collected into one large table using data.table's 'rbindlist'. \$metadata is a little weirder because the metadata doesn't fit neatly into a tidy format but things are hopefully named helpfully. \$chroms was added in v1.3 and contains 7 columns: chromatogram type (usually

TIC, BPC or SRM info), chromatogram index, target mz, product mz, retention time (rt), and intensity (int). \$DAD was also added in v1.3 and contains has three columns: retention time (rt), wavelength (lambda), and intensity (int). Data requested that does not exist in the provided files (such as MS2 data requested from MS1-only files) will return an empty (zero-row) data.table.

Examples

```
library(RaMS)
# Extract MS1 data from a couple files
sample_dir <- system.file("extdata", package = "RaMS")
sample_files <- list.files(sample_dir, full.names=TRUE)
multifile_data <- grabMSdata(sample_files[c(3, 5, 6)], grab_what="MS1")

# "Stream" data from the internet (i.e. Metabolights)
## Not run:
access_url <- "https://www.ebi.ac.uk/metabolights/MTBLS703/files"

# URL below obtained by right-clicking site download button and copying
# link address
sample_url <- paste0("https://www.ebi.ac.uk/metabolights/ws/studies/",
                    "MTBLS703/download/acefcd61-a634-4f35-9c3c-c572",
                    "ade5acf3?file=161024_Smp_LB12HL_AB_pos.mzXML")
file_data <- grabMSdata(sample_url, grab_what="everything", verbosity=2)

## End(Not run)
```

grabMzmlBPC

Grab the BPC or TIC from a file

Description

The base peak intensity and total ion current are actually written into the mzML files and aren't encoded, making retrieval of BPC and TIC information blazingly fast if parsed correctly.

Usage

```
grabMzmlBPC(xml_data, rrange, TIC = FALSE)
```

Arguments

xml_data	An 'xml2' nodeset, usually created by applying 'read_xml' to an mzML file.
rrange	A vector of length 2 containing an upper and lower bound on retention times of interest. Providing a range here can speed up load times (although not enormously, as the entire file must still be read) and reduce the final object's size.
TIC	Boolean. If TRUE, the TIC is extracted rather than the BPC.

Value

A 'data.table' with columns for retention time (rt), and intensity (int).

grabMzmlDAD	<i>Extract the DAD data from an mzML nodeset</i>
-------------	--

Description

Extract the DAD data from an mzML nodeset

Usage

```
grabMzmlDAD(xml_data, rtrange, file_metadata)
```

Arguments

xml_data	An 'xml2' nodeset, usually created by applying 'read_xml' to an mzML file.
rtrange	A vector of length 2 containing an upper and lower bound on retention times of interest. Providing a range here can speed up load times (although not enormously, as the entire file must still be read) and reduce the final object's size.
file_metadata	Information about the file used to decode the binary arrays containing m/z and intensity information.

Value

A 'data.table' with columns for retention time (rt), wavelength (lambda), and intensity (int).

grabMzmlData	<i>Get mass-spectrometry data from an mzML file</i>
--------------	---

Description

This function handles the mzML side of things, reading in files that are written in the mzML format. Much of the code is similar to the mzXML format, but the xpath handles are different and the mz/int array is encoded as two separate entries rather than simultaneously. This function has been exposed to the user in case per-file optimization (such as peakpicking or additional filtering) is desired before the full data object is returned.

Usage

```
grabMzmlData(
  filename,
  grab_what,
  verbosity = 0,
  mz = NULL,
  ppm = NULL,
  rtrange = NULL,
  prefilter = -1
)
```


Arguments

filename	A single filename to read into R's memory. Both absolute and relative paths are acceptable.
grab_what	What data should be read from the file? Options include "MS1" for data only from the first spectrometer, "MS2" for fragmentation data, "BPC" for rapid access to the base peak chromatogram, "TIC" for rapid access to the total ion chromatogram, "DAD" for DAD (UV) data, and "chroms" for precompiled chromatogram data (especially useful for MRM but often contains BPC/TIC in other files). Metadata can be accessed with "metadata", which provides information about the instrument and time the file was run. These options can be combined (i.e. 'grab_data=c("MS1", "MS2", "BPC")') or this argument can be set to "everything" to extract all of the above. Options "EIC" and "EIC_MS2" are useful when working with files whose total size exceeds working memory - it first extracts all relevant MS1 and MS2 data, respectively, then discards data outside of the mass range(s) calculated from the provided mz and ppm. The default, "everything", includes all MS1, MS2, BPC, TIC, and metadata.
verbosity	Three levels of processing output to the R console are available, with increasing verbosity corresponding to higher integers. A verbosity of zero means that no output will be produced, useful when wrapping within larger functions. A verbosity of 1 will produce a progress bar using base R's txtProgressBar function. A verbosity of 2 or higher will produce timing output for each individual file read in.
mz	A vector of the mass-to-charge ratio for compounds of interest. Only used when combined with 'grab_what = "EIC"' (see above). Multiple masses can be provided.
ppm	A single number corresponding to the mass accuracy (in parts per million) of the instrument on which the data was collected. Only used when combined with 'grab_what = "EIC"' (see above).
rtrange	A vector of length 2 containing an upper and lower bound on retention times of interest. Providing a range here can speed up load times (although not enormously, as the entire file must still be read) and reduce the final object's size.
prefilter	A single number corresponding to the minimum intensity of interest in the MS1 data. Data points with intensities below this threshold will be silently dropped, which can dramatically reduce the size of the final object. Currently only works with MS1 data, but could be expanded easily to handle more.

Value

A list of 'data.table's, each named after the arguments requested in grab_what. E.g. \$MS1 contains MS1 information, \$MS2 contains fragmentation info, etc. MS1 data has four columns: retention time (rt), mass-to-charge (mz), intensity (int), and filename. MS2 data has six: retention time (rt), precursor m/z (prezmz), fragment m/z (fragmz), fragment intensity (int), collision energy (voltage), and filename. Data requested that does not exist in the provided files (such as MS2 data requested from MS1-only files) will return an empty (length zero) data.table. The data.tables extracted from each of the individual files are collected into one large table using data.table's 'rbindlist'. \$metadata is a little weirder because the metadata doesn't fit neatly into a tidy format but things are hopefully

named helpfully. \$Schroms was added in v1.3 and contains 7 columns: chromatogram type (usually TIC, BPC or SRM info), chromatogram index, target mz, product mz, retention time (rt), and intensity (int). \$DAD was also added in v1.3 and contains three columns: retention time (rt), wavelength (lambda), and intensity (int). Data requested that does not exist in the provided files (such as MS2 data requested from MS1-only files) will return an empty (zero-row) data.table.

Examples

```
sample_file <- system.file("extdata", "LB12HL_AB.mzML.gz", package = "RaMS")
file_data <- grabMzmlData(sample_file, grab_what="MS1")
## Not run:
# Extract MS1 data and a base peak chromatogram
file_data <- grabMzmlData(sample_file, grab_what=c("MS1", "BPC"))
# Extract data from a retention time subset
file_data <- grabMzmlData(sample_file, grab_what=c("MS1", "BPC"),
                          rtrange=c(5, 7))
# Extract EIC for a specific mass
file_data <- grabMzmlData(sample_file, grab_what="EIC", mz=118.0865, ppm=5)
# Extract EIC for several masses simultaneously
file_data <- grabMzmlData(sample_file, grab_what="EIC", ppm=5,
                          mz=c(118.0865, 146.118104, 189.123918))

# Extract MS2 data
sample_file <- system.file("extdata", "DDApos_2.mzML.gz", package = "RaMS")
MS2_data <- grabMzmlData(sample_file, grab_what="MS2")

## End(Not run)
```

grabMzmlEncodingData *Helper function to extract mzML file encoding data*

Description

Helper function to extract mzML file encoding data

Usage

```
grabMzmlEncodingData(xml_data)
```

Arguments

xml_data mzML data as parsed by xml2

Value

A list of values used by other parsing functions, currently compression, mz_precision, int_precision

grabMzmlMetadata	<i>Helper function to extract mzML file metadata</i>
------------------	--

Description

Helper function to extract mzML file metadata

Usage

```
grabMzmlMetadata(xml_data)
```

Arguments

xml_data	mzML data as parsed by xml2
----------	-----------------------------

Value

A list of values corresponding to various pieces of metadata for each file

grabMzmlMS1	<i>Extract the MS1 data from an mzML nodeset</i>
-------------	--

Description

Extract the MS1 data from an mzML nodeset

Usage

```
grabMzmlMS1(xml_data, rtrange, file_metadata, prefilter)
```

Arguments

xml_data	An 'xml2' nodeset, usually created by applying 'read_xml' to an mzML file.
rtrange	A vector of length 2 containing an upper and lower bound on retention times of interest. Providing a range here can speed up load times (although not enormously, as the entire file must still be read) and reduce the final object's size.
file_metadata	Information about the file used to decode the binary arrays containing m/z and intensity information.
prefilter	The lowest intensity value of interest, used to reduce file size (and especially useful for profile mode data with many 0 values)

Value

A 'data.table' with columns for retention time (rt), m/z (mz), and intensity (int).

 grabMzmlMS2

Extract the MS2 data from an mzML nodeset

Description

Extract the MS2 data from an mzML nodeset

Usage

```
grabMzmlMS2(xml_data, rrange, file_metadata)
```

Arguments

xml_data	An ‘xml2’ nodeset, usually created by applying ‘read_xml’ to an mzML file.
rrange	A vector of length 2 containing an upper and lower bound on retention times of interest. Providing a range here can speed up load times (although not enormously, as the entire file must still be read) and reduce the final object’s size.
file_metadata	Information about the file used to decode the binary arrays containing m/z and intensity information.

Value

A ‘data.table’ with columns for retention time (rt), precursor m/z (mz), fragment m/z (fragmz), collision energy (voltage), and intensity (int).

 grabMzxmlBPC

Grab the BPC or TIC from a file

Description

The base peak intensity and total ion current are actually written into the mzXML files and aren’t encoded, making retrieval of BPC and TIC information blazingly fast if parsed correctly.

Usage

```
grabMzxmlBPC(xml_data, TIC = FALSE, rrange)
```

Arguments

xml_data	An ‘xml2’ nodeset, usually created by applying ‘read_xml’ to an mzML file.
TIC	Boolean. If TRUE, the TIC is extracted rather than the BPC.
rrange	A vector of length 2 containing an upper and lower bound on retention times of interest. Providing a range here can speed up load times (although not enormously, as the entire file must still be read) and reduce the final object’s size.

Value

A 'data.table' with columns for retention time (rt), and intensity (int).

grabMzxmlData	<i>Get mass-spectrometry data from an mzXML file</i>
---------------	--

Description

This function handles the mzXML side of things, reading in files that are written in the mzXML format. Much of the code is similar to the mzXML format, but the xpath handles are different and the mz/int array is encoded simultaneously rather than as two separate entries. This function has been exposed to the user in case per-file optimization (such as peakpicking or additional filtering) is desired before the full data object is returned.

Usage

```
grabMzxmlData(
  filename,
  grab_what,
  verbosity = 0,
  rtrange = NULL,
  mz = NULL,
  ppm = NULL,
  prefilter = -1
)
```

Arguments

filename	A single filename to read into R's memory. Both absolute and relative paths are acceptable.
grab_what	What data should be read from the file? Options include "MS1" for data only from the first spectrometer, "MS2" for fragmentation data, "BPC" for rapid access to the base peak chromatogram, and "TIC" for rapid access to the total ion chromatogram. DAD and chromatogram ("DAD" and "chroms") are unavailable for mzXML files. Metadata can be accessed with "metadata", which provides information about the instrument and time the file was run. These options can be combined (i.e. 'grab_data=c("MS1", "MS2", "BPC")') or this argument can be set to "everything" to extract all of the above. Options "EIC" and "EIC_MS2" are useful when working with files whose total size exceeds working memory - it first extracts all relevant MS1 and MS2 data, respectively, then discards data outside of the mass range(s) calculated from the provided mz and ppm. The default, "everything", includes all MS1, MS2, BPC, TIC, and metadata.
verbosity	Three levels of processing output to the R console are available, with increasing verbosity corresponding to higher integers. A verbosity of zero means that no output will be produced, useful when wrapping within larger functions. A verbosity of 1 will produce a progress bar using base R's txtProgressBar function.

	A verbosity of 2 or higher will produce timing output for each individual file read in.
rtrange	Not supported for mzXML data. Only provided here so as to throw a friendly warning rather than an unexpected error.
mz	A vector of the mass-to-charge ratio for compounds of interest. Only used when combined with 'grab_what = "EIC"' (see above). Multiple masses can be provided.
ppm	A single number corresponding to the mass accuracy (in parts per million) of the instrument on which the data was collected. Only used when combined with 'grab_what = "EIC"' (see above).
prefilter	A single number corresponding to the minimum intensity of interest in the MS1 data. Data points with intensities below this threshold will be silently dropped, which can dramatically reduce the size of the final object. Currently only works with MS1 data, but could be expanded easily to handle more.

Value

A list of 'data.table's, each named after the arguments requested in grab_what. E.g. \$MS1 contains MS1 information, \$MS2 contains fragmentation info, etc. MS1 data has four columns: retention time (rt), mass-to-charge (mz), intensity (int), and filename. MS2 data has six: retention time (rt), precursor m/z (prezmz), fragment m/z (fragmz), fragment intensity (int), collision energy (voltage), and filename. Data requested that does not exist in the provided files (such as MS2 data requested from MS1-only files) will return an empty (length zero) data.table. The data.tables extracted from each of the individual files are collected into one large table using data.table's 'rbindlist'. \$metadata is a little weirder because the metadata doesn't fit neatly into a tidy format but things are hopefully named helpfully. Data requested that does not exist in the provided files (such as DAD or chromatogram data) will return an empty (zero-row) data.table.

Examples

```
sample_file <- system.file("extdata", "LB12HL_AB.mzXML.gz", package = "RaMS")
file_data <- grabMzxmlData(sample_file, grab_what="MS1")
## Not run:
# Extract MS1 data and a base peak chromatogram
file_data <- grabMzxmlData(sample_file, grab_what=c("MS1", "BPC"))
# Extract EIC for a specific mass
file_data <- grabMzxmlData(sample_file, grab_what="EIC", mz=118.0865, ppm=5)
# Extract EIC for several masses simultaneously
file_data <- grabMzxmlData(sample_file, grab_what="EIC", ppm=5,
                           mz=c(118.0865, 146.118104, 189.123918))

# Extract MS2 data
sample_file <- system.file("extdata", "DDApos_2.mzXML.gz", package = "RaMS")
MS2_data <- grabMzxmlData(sample_file, grab_what="MS2")

## End(Not run)
```

grabMzxmlEncodingData *Helper function to extract mzXML file metadata*

Description

Helper function to extract mzXML file metadata

Usage

```
grabMzxmlEncodingData(xml_data)
```

Arguments

xml_data mzXML data as parsed by xml2

Value

A list of values used by other parsing functions, currently compression, precision, and endian encoding (endi_enc)

grabMzxmlMetadata *Helper function to extract mzXML file metadata*

Description

Helper function to extract mzXML file metadata

Usage

```
grabMzxmlMetadata(xml_data)
```

Arguments

xml_data mzXML data as parsed by xml2

Value

A list of values corresponding to various pieces of metadata for each file

grabMzxmlMS1	<i>Extract the MS1 data from an mzXML nodeset</i>
--------------	---

Description

Extract the MS1 data from an mzXML nodeset

Usage

```
grabMzxmlMS1(xml_data, file_metadata, rrange, prefilter)
```

Arguments

xml_data	An 'xml2' nodeset, usually created by applying 'read_xml' to an mzXML file.
file_metadata	Information about the file used to decode the binary arrays containing m/z and intensity information.
rrange	A vector of length 2 containing an upper and lower bound on retention times of interest. Providing a range here can speed up load times (although not enormously, as the entire file must still be read) and reduce the final object's size.
prefilter	The lowest intensity value of interest, used to reduce file size (and especially useful for profile mode data with many 0 values)

Value

A 'data.table' with columns for retention time (rt), m/z (mz), and intensity (int).

grabMzxmlMS2	<i>Extract the MS2 data from an mzXML nodeset</i>
--------------	---

Description

Extract the MS2 data from an mzXML nodeset

Usage

```
grabMzxmlMS2(xml_data, file_metadata, rrange)
```

Arguments

xml_data	An 'xml2' nodeset, usually created by applying 'read_xml' to an mzXML file.
file_metadata	Information about the file used to decode the binary arrays containing m/z and intensity information.
rrange	A vector of length 2 containing an upper and lower bound on retention times of interest. Providing a range here can speed up load times (although not enormously, as the entire file must still be read) and reduce the final object's size.

Value

A 'data.table' with columns for retention time (rt), precursor m/z (mz), fragment m/z (fragmz), collision energy (voltage), and intensity (int).

grabMzxmlSpectraMzInt *Extract the mass-to-charge data from the spectra of an mzXML node-set*

Description

The mz and intensity information of mzXML files are encoded as a binary array, sometimes compressed via gzip or zlib or numpress. This code finds all the m/z-int binary arrays and converts them back to the original measurements. See <https://github.com/ProteoWizard/pwiz/issues/1301>

Usage

```
grabMzxmlSpectraMzInt(xml_nodes, file_metadata)
```

Arguments

xml_nodes	An xml_nodeset object corresponding to the spectra collected by the mass spectrometer, usually produced by applying 'xml_find_all' to an MS1 or MS2 node-set.
file_metadata	Information about the file used to decode the binary arrays containing m/z and intensity information. Here, the compression and mz precision information is relevant.

Value

A numeric vector of masses, many for each scan.

grabMzxmlSpectraPremz *Extract the precursor mass from the spectra of an mzXML nodeset*

Description

Extract the precursor mass from the spectra of an mzXML nodeset

Usage

```
grabMzxmlSpectraPremz(xml_nodes)
```

Arguments

xml_nodes	An xml_nodeset object corresponding to the spectra collected by the mass spectrometer, usually produced by applying 'xml_find_all' to an MS1 or MS2 node-set.
-----------	---

Value

A numeric vector of precursor masses, one for each scan

grabMzxmlSpectraRt *Extract the retention time from the spectra of an mzXML nodeset*

Description

Extract the retention time from the spectra of an mzXML nodeset

Usage

```
grabMzxmlSpectraRt(xml_nodes)
```

Arguments

xml_nodes An xml_nodeset object corresponding to the spectra collected by the mass spectrometer, usually produced by applying 'xml_find_all' to an MS1 or MS2 nodeset.

Value

A numeric vector of retention times, one for each scan

grabMzxmlSpectraVoltage
 Extract the collision energies from the spectra of an mzXML nodeset

Description

Although the collision energy is typically fixed per file, it's equally fast (afaik) to just grab them all individually here. Also, I'm worried about these rumors of "ramped" collision energies

Usage

```
grabMzxmlSpectraVoltage(xml_nodes)
```

Arguments

xml_nodes An xml_nodeset object corresponding to the spectra collected by the mass spectrometer, usually produced by applying 'xml_find_all' to an MS1 or MS2 nodeset.

Value

A numeric vector of collision energies, one for each scan.

grabSpectraInt	<i>Extract the intensity information from the spectra of an mzML nodeset</i>
----------------	--

Description

The m/z and intensity information of mzML files are encoded as binary arrays, sometimes compressed via gzip or zlib or numpress. This code finds all the intensity binary arrays and converts them back to the original measurements. See <https://github.com/ProteoWizard/pwiz/issues/1301>

Usage

```
grabSpectraInt(xml_nodes, file_metadata)
```

Arguments

xml_nodes	An xml_nodeset object corresponding to the spectra collected by the mass spectrometer, usually produced by applying 'xml_find_all' to an MS1 or MS2 nodeset.
file_metadata	Information about the file used to decode the binary arrays containing m/z and intensity information. Here, the compression and int precision information is relevant.

Value

A numeric vector of intensities, many for each scan.

grabSpectraMz	<i>Extract the mass-to-charge data from the spectra of an mzML nodeset</i>
---------------	--

Description

The m/z and intensity information of mzML files are encoded as binary arrays, sometimes compressed via gzip or zlib or numpress. This code finds all the m/z binary arrays and converts them back to the original measurements. See <https://github.com/ProteoWizard/pwiz/issues/1301>

Usage

```
grabSpectraMz(xml_nodes, file_metadata)
```

Arguments

xml_nodes	An xml_nodeset object corresponding to the spectra collected by the mass spectrometer, usually produced by applying 'xml_find_all' to an MS1 or MS2 nodeset.
file_metadata	Information about the file used to decode the binary arrays containing m/z and intensity information. Here, the compression and m/z precision information is relevant.

Value

A numeric vector of masses, many for each scan.

grabSpectraPremz	<i>Extract the precursor mass from the spectra of an mzML nodeset</i>
------------------	---

Description

Extract the precursor mass from the spectra of an mzML nodeset

Usage

```
grabSpectraPremz(xml_nodes)
```

Arguments

xml_nodes	An xml_nodeset object corresponding to the spectra collected by the mass spectrometer, usually produced by applying 'xml_find_all' to an MS1 or MS2 nodeset.
-----------	--

Value

A numeric vector of precursor masses, one for each scan

grabSpectraRt	<i>Extract the retention time from the spectra of an mzML nodeset</i>
---------------	---

Description

Extract the retention time from the spectra of an mzML nodeset

Usage

```
grabSpectraRt(xml_nodes)
```

Arguments

xml_nodes	An xml_nodeset object corresponding to the spectra collected by the mass spectrometer, usually produced by applying 'xml_find_all' to an MS1 or MS2 nodeset.
-----------	--

Value

A numeric vector of retention times, one for each scan

grabSpectraVoltage	<i>Extract the collision energies from the spectra of an mzML nodeset</i>
--------------------	---

Description

Although the collision energy is typically fixed per file, it's equally fast (afaik) to just grab them all individually here. Also, I'm worried about these rumors of "ramped" collision energies

Usage

```
grabSpectraVoltage(xml_nodes)
```

Arguments

xml_nodes	An xml_nodeset object corresponding to the spectra collected by the mass spectrometer, usually produced by applying 'xml_find_all' to an MS1 or MS2 nodeset.
-----------	--

Value

A numeric vector of collision energies, one for each scan.

minifyMSdata	<i>Shrink MS data by including only data points near masses of interest</i>
--------------	---

Description

MS files can be annoyingly large if only a few masses are of interest. This large size makes it difficult to share them online for debugging purposes and often means that untargeted algorithms spend a lot of time picking peaks in data that's irrelevant. minifyMSdata is a function designed to "minify" MS files by extracting only those data points that are within the ppm error of an m/z value of interest, and returns the file essentially otherwise unchanged.

Usage

```
minifyMSdata(  
  files,  
  output_files = NULL,  
  mz_exclude = NULL,  
  mz_include = NULL,  
  ppm = NULL,  
  warn = TRUE,  
  prefilter = -1,  
  verbosity = NULL  
)
```

Arguments

files	The name of a single file to be minified, usually produced by Proteowizard's 'msconvert' or something similar.
output_files	The name of the file to be written out.
mz_exclude	A vector of m/z values that should be excluded from the minified file. This argument must be used with the 'ppm' argument and should not be used with mz_include. For each mass provided, an m/z window of +/- 'ppm' is calculated, and all data points within that window are removed.
mz_include	A vector of m/z values that should be included in the minified file. This argument must be used with the 'ppm' argument and should not be used with mz_exclude. For each mass provided, an m/z window of +/- 'ppm' is calculated, and all data points within that window are kept.
ppm	The parts-per-million error of the instrument used to collect the original file.
warn	Boolean. Should the function warn the user when removing an index from an mzML file?
prefilter	A single number corresponding to the minimum intensity of interest in the MS1 data. Data points with intensities below this threshold will be silently dropped, which can dramatically reduce the size of the final object. Currently only works with MS1 data, but could be expanded easily to handle more.
verbosity	A single number with a sensible default behavior. If larger than 2, will render a progress bar as files are processed.

Value

Invisibly, the name of the new files.

Examples

```
## Not run:
library(RaMS)
# Extract data corresponding to only valine and homarine
# m/z = 118.0865 and 138.0555, respectively
filename <- system.file("extdata", "LB12HL_AB.mzML.gz", package = "RaMS")
output_filename <- "mini_LB12HL_AB.mzML"
include_mzs <- c(118.0865, 138.0555)
minifyMSdata(filename, output_filename, mz_include=include_mzs, ppm=5)
unlink(output_filename)

# Exclude data corresponding to valine and homarine
filename <- system.file("extdata", "LB12HL_AB.mzML.gz", package = "RaMS")
output_filename <- "mini_LB12HL_AB.mzML"
exclude_mzs <- c(118.0865, 138.0555)
minifyMSdata(filename, output_filename, mz_exclude=exclude_mzs, ppm=5)
unlink(output_filename)

## End(Not run)
```

`minifyMzml`*Shrink mzML files by including only data points near masses of interest*

Description

mzML files can be annoyingly large if only a few masses are of interest. This large size makes it difficult to share them online for debugging purposes and often means that untargeted algorithms spend a lot of time picking peaks in data that's irrelevant. `minifyMzml` is a function designed to "minify" mzML files by extracting only those data points that are within a ppm error of an m/z value of interest, and returns the file essentially otherwise unchanged. This function currently works only on MS1 data, but is reasonably expandable if demand becomes evident.

Usage

```
minifyMzml(  
  filename,  
  output_filename,  
  ppm,  
  mz_exclude = NULL,  
  mz_include = NULL,  
  warn = TRUE,  
  prefilter = -1  
)
```

Arguments

<code>filename</code>	The name of a single file to be minified, usually produced by Proteowizard's 'msconvert' or something similar.
<code>output_filename</code>	The name of the file to be written out.
<code>ppm</code>	The parts-per-million error of the instrument used to collect the original file.
<code>mz_exclude</code>	A vector of m/z values that should be excluded from the minified file. This argument must be used with the 'ppm' argument and should not be used with <code>mz_include</code> . For each mass provided, an m/z window of +/- 'ppm' is calculated, and all data points within that window are removed.
<code>mz_include</code>	A vector of m/z values that should be included in the minified file. This argument must be used with the 'ppm' argument and should not be used with <code>mz_exclude</code> . For each mass provided, an m/z window of +/- 'ppm' is calculated, and all data points within that window are kept.
<code>warn</code>	Boolean. Should the function warn the user when removing an index from an mzML file?
<code>prefilter</code>	A single number corresponding to the minimum intensity of interest in the MS1 data. Data points with intensities below this threshold will be silently dropped, which can dramatically reduce the size of the final object. Currently only works with MS1 data, but could be expanded easily to handle more.

Value

Invisibly, the name of the new file.

Examples

```
## Not run:
library(RaMS)
# Extract data corresponding to only valine and homarine
# m/z = 118.0865 and 138.0555, respectively
filename <- system.file("extdata", "LB12HL_AB.mzML.gz", package = "RaMS")
output_filename <- "mini_LB12HL_AB.mzML"
include_mzs <- c(118.0865, 138.0555)
minifyMzml(filename, output_filename, mz_include=include_mzs, ppm=5)
unlink(output_filename)

# Exclude data corresponding to valine and homarine
filename <- system.file("extdata", "LB12HL_AB.mzML.gz", package = "RaMS")
output_filename <- "mini_LB12HL_AB.mzML"
exclude_mzs <- c(118.0865, 138.0555)
minifyMzml(filename, output_filename, mz_exclude=exclude_mzs, ppm=5)
unlink(output_filename)

## End(Not run)
```

minifyMzxml

Shrink mzXML files by including only data points near masses of interest

Description

mzXML files can be annoyingly large if only a few masses are of interest. This large size makes it difficult to share them online for debugging purposes and often means that untargeted algorithms spend a lot of time picking peaks in data that's irrelevant. `minifyMzxml` is a function designed to "minify" mzXML files by extracting only those data points that are within a ppm error of an m/z value of interest, and returns the file essentially otherwise unchanged. This function currently works only on MS1 data, but is reasonably expandable if demand becomes evident.

Usage

```
minifyMzxml(
  filename,
  output_filename,
  ppm,
  mz_exclude = NULL,
  mz_include = NULL,
  prefilter = -1,
  warn = TRUE
)
```


Arguments

filename	The name of a single file to be minified, usually produced by Proteowizard's 'msconvert' or something similar.
output_filename	The name of the file to be written out.
ppm	The parts-per-million error of the instrument used to collect the original file.
mz_exclude	A vector of m/z values that should be excluded from the minified file. This argument must be used with the 'ppm' argument and should not be used with mz_include. For each mass provided, an m/z window of +/- 'ppm' is calculated, and all data points within that window are removed.
mz_include	A vector of m/z values that should be included in the minified file. This argument must be used with the 'ppm' argument and should not be used with mz_exclude. For each mass provided, an m/z window of +/- 'ppm' is calculated, and all data points within that window are kept.
prefilter	A single number corresponding to the minimum intensity of interest in the MS1 data. Data points with intensities below this threshold will be silently dropped, which can dramatically reduce the size of the final object. Currently only works with MS1 data, but could be expanded easily to handle more.
warn	Boolean. Should the function warn the user when removing an index from an mzML file?

Value

Invisibly, the name of the new file.

Examples

```
## Not run:
library(RaMS)
# Extract data corresponding to only valine and homarine
# m/z = 118.0865 and 138.0555, respectively
filename <- system.file("extdata", "LB12HL_AB.mzXML.gz", package = "RaMS")
output_filename <- "mini_LB12HL_AB.mzXML"
include_mzs <- c(118.0865, 138.0555)
minifyMzxml(filename, output_filename, mz_include=include_mzs, ppm=5)
unlink(output_filename)

# Exclude data corresponding to valine and homarine
filename <- system.file("extdata", "LB12HL_AB.mzXML.gz", package = "RaMS")
output_filename <- "mini_LB12HL_AB.mzXML"
exclude_mzs <- c(118.0865, 138.0555)
minifyMzxml(filename, output_filename, mz_exclude=exclude_mzs, ppm=5)
unlink(output_filename)

## End(Not run)
```

msdata_connection	<i>S3 constructor for msdata_connection</i>
-------------------	---

Description

S3 constructor for msdata_connection

Usage

```
msdata_connection(x)
```

Arguments

x	This is a thing?
---	------------------

Value

Itself, with the class?

node2dt	<i>Convert node to data.table</i>
---------	-----------------------------------

Description

Convert node to data.table

Usage

```
node2dt(dubset_node, ms_level)
```

Arguments

dubset_node	The "data subset" node with children rt, mz, etc.
ms_level	The requested MS level to search for

Value

A data.table with columns depending on the MS level requested

pmpm *Plus/minus parts per million*

Description

It shouldn't be hard to translate a point mass into a mass window bounded by spectrometer accuracy.

Usage

```
pmpm(mass, ppm = 4)
```

Arguments

mass	A length-1 numeric representing the mass of interest for which a mass range is desired.
ppm	The parts-per-million accuracy of the mass spectrometer on which the data was collected.

Value

A length-2 numeric representing the mass range requested

Examples

```
pmpm(100, 5)
pmpm(1000000, 5)
pmpm(118.0865, 2.5)
pmpm(892.535313, 10)
```

```
print.msdata_connection
S3 print option for msdata_connection objects
```

Description

S3 print option for msdata_connection objects

Usage

```
## S3 method for class 'msdata_connection'
print(x, ...)
```

Arguments

x	An msdata_connection object containing files and grab_what
...	Other arguments to be passed to print.default, I guess

Value

Messages, mostly

tmzmlMaker

Maker of tmzML documents

Description

This function converts mzML and mzXML documents into "transposed" mzML (tmzML) documents. Traditional mass-spec data is organized by scan number, corresponding to retention time, but this isn't always the most sensible format. Often, it makes more sense to organize a mass-spec file by m/z ratio instead. This allows parsers to scan and decode a much smaller portion of the file when searching for a specific mass, as opposed to the traditional format which requires that every scan be opened, searched, and subset. The tmzML document implements this strategy and allows the creation of MS object representations that use essentially zero memory because the data is read off the disk instead of being stored in RAM. RaMS has been designed to interface with these new file types identically to traditional files, allowing all your favorite tidyverse tricks to work just as well and much more quickly.

Usage

```
tmzmlMaker(input_filename, output_filename = NULL, verbosity = 0, binwidth = 3)
```

Arguments

- `input_filename` Character vector of length 1 with the name of the file to be converted. Can only handle mzML and mzXML currently - other formats should be converted to one of these first, using (for example) Proteowizard's msconvert tool.
- `output_filename` The name of the file that will be written out. Should end in ".tmzML" and will throw a warning otherwise. Often, it makes sense to have two folders in a working directory, one containing the original mzML files and a second, parallel folder for the tmzMLs.
- `verbosity` Numeric value between 0 and 2, corresponding to level of verbosity shared by the function as it proceeds. 0 means no output, 1 will produce mile markers after file opening, MS1 and MS2 conversion, and 2 will provide progress bars between each mile marker.
- `binwidth` Numeric value controlling the width of the bins in m/z space to create. Because MS data is created in such a way that m/z values are continuous, they must be binned together to create a discrete representation that can be searched efficiently. Lower values (0.1-1) will have faster retrieval times, while higher values (5-10) will have faster conversion times.

Value

An `msdata_connection` object. This object behaves exactly like a normal RaMS list with values for MS1, MS2, etc. but secretly just contains pointers to the files requested because the data is extracted on the fly. The S3 `msdata_connection` object is necessary to create new behaviors for '\$' and '[' that allow indexing like normal.

Examples

```
## Not run:
sample_dir <- system.file("extdata", package = "RaMS")
sample_files <- list.files(sample_dir, full.names=TRUE, pattern="LB.*mzML")
tmzml_filenames <- gsub(x=sample_files, "\\.*mzML.gz", ".tmzML")

# Convert a single file
tmzmlMaker(sample_files[1], tmzml_filenames[1])
file_data <- grabMSdata(tmzml_filenames[1], grab_what="everything", verbosity=2)
file_data$MS1[mz%between%pmpm(118.0865)]

# Multiple files
mapply(tmzmlMaker, sample_files, tmzml_filenames)
file_data <- grabMSdata(tmzml_filenames, grab_what="everything", verbosity=2)
betaine_data <- file_data$MS1[mz%between%pmpm(118.0865)]

# Plot output
plot(betaine_data$rt, betaine_data$int, type="l")
library(ggplot2)
ggplot(betaine_data) + geom_line(aes(x=rt, y=int, color=filename))

# Clean up afterward
file.remove(tmzml_filenames)

## End(Not run)
```

[.msdata_connection *S3 indexing for msdata_connection objects*

Description

This is the step that actually performs the file opening and extraction!

Usage

```
## S3 method for class 'msdata_connection'
msdata_obj[sub_func]
```

Arguments

<code>msdata_obj</code>	An <code>msdata_connection</code> object containing files and <code>grab_what</code>
<code>sub_func</code>	The function that will be parsed and used to subset the file

Value

A data.table with columns rt, mz, int, and filename

\$.msdata_connection *S3 dollar sign notation for msdata_connection objects*

Description

S3 dollar sign notation for msdata_connection objects

Usage

```
## S3 method for class 'msdata_connection'  
msdata_obj$ms_level
```

Arguments

msdata_obj An msdata_connection object containing files and grab_what
ms_level The requested MS level of the object

Value

An msdata_connection object with only a single MS level

Index

[.msdata_connection, 29
\$.msdata_connection, 30

checkOutputQuality, 2

getEncoded, 3
giveEncoding, 4
grabAccessionData, 4
grabMSdata, 5
grabMzmlBPC, 7
grabMzmlDAD, 8
grabMzmlData, 8
grabMzmlEncodingData, 10
grabMzmlMetadata, 11
grabMzmlMS1, 11
grabMzmlMS2, 12
grabMzxmlBPC, 12
grabMzxmlData, 13
grabMzxmlEncodingData, 15
grabMzxmlMetadata, 15
grabMzxmlMS1, 16
grabMzxmlMS2, 16
grabMzxmlSpectraMzInt, 17
grabMzxmlSpectraPremz, 17
grabMzxmlSpectraRt, 18
grabMzxmlSpectraVoltage, 18
grabSpectraInt, 19
grabSpectraMz, 19
grabSpectraPremz, 20
grabSpectraRt, 20
grabSpectraVoltage, 21

minifyMSdata, 21
minifyMzml, 23
minifyMzxml, 24
msdata_connection, 26

node2dt, 26

pmpm, 27
print.msdata_connection, 27

tmzmlMaker, 28