

# Package ‘chem16S’

July 17, 2023

**Date** 2023-07-17

**Version** 1.0.0

**Title** Chemical Metrics of Community Reference Proteomes

**Author** Jeffrey Dick [aut, cre] (<<https://orcid.org/0000-0002-0687-5890>>)

**Maintainer** Jeffrey Dick <j3ffdic@gmail.com>

**Depends** R (>= 3.5.0)

**Imports** plyr, ggplot2, rlang, reshape2, phyloseq

**Suggests** tinytest, knitr, patchwork, ggpmisc, rmarkdown

**Description** Combines taxonomic classifications of high-throughput 16S rRNA gene sequences with reference proteomes of archaeal and bacterial taxa to generate amino acid compositions of community reference proteomes. Calculates chemical metrics including carbon oxidation state ('Zc'), stoichiometric oxidation and hydration state ('nO2' and 'nH2O'), H/C, N/C, O/C, and S/C ratios, grand average of hydropathicity ('GRAVY'), isoelectric point ('pI'), protein length, and average molecular weight of amino acid residues. Uses precomputed reference proteomes for archaea and bacteria derived from the Genome Taxonomy Database ('GTDB'). Also includes reference proteomes derived from the NCBI Reference Sequence ('RefSeq') database and manual mapping from the 'RDP Classifier' training set to 'RefSeq' taxonomy as described by Dick and Tan (2023) <[doi:10.1007/s00248-022-01988-9](https://doi.org/10.1007/s00248-022-01988-9)>. Processes taxonomic classifications in 'RDP Classifier' format or OTU tables in 'phyloseq-class' objects from the Bioconductor package 'phyloseq'.

**Encoding** UTF-8

**License** GPL-3

**VignetteBuilder** knitr

**URL** <https://github.com/jedick/chem16S>,  
<https://doi.org/10.5281/zenodo.6793059>

**NeedsCompilation** no

**Repository** CRAN

**Date/Publication** 2023-07-17 17:10:02 UTC

## R topics documented:

chem16S-package . . . . .	2
calc_metrics . . . . .	5
chem16S-data . . . . .	7
chemlab . . . . .	8
get_metadata . . . . .	9
get_metrics . . . . .	10
map_taxa . . . . .	12
physeq . . . . .	14
plot_metrics . . . . .	16
read_RDP . . . . .	18

<b>Index</b>	<b>20</b>
--------------	-----------

---

chem16S-package	<i>Chemical metrics of community reference proteomes</i>
-----------------	--

---

### Description

Functions and data to calculate chemical metrics for reference proteomes for microbial (archaeal and bacterial) communities. Amino acid compositions of community reference proteomes are generated by combining reference proteomes of taxa (derived from GTDB or RefSeq) with taxonomic classifications of 16S rRNA gene sequences.

### Details

- `read_RDP` - Read and filter an RDP Classifier table
- `map_taxa` - Map RDP Classifier assignments to the NCBI taxonomy
- `get_metrics` - Get chemical metrics for community reference proteomes
- `get_metadata` - Example function for retrieving sample metadata
- `plot_metrics` - Plot metrics with symbols and colors based on metadata
- `physeq` - Functions designed to analyze `phyloseq-class` objects

#### History:

1. Work begun in 2021. The combination of RefSeq reference proteomes with taxonomic abundances to compute community-level chemical metrics was described by Dick and Tan (2023). **chem16S** originated as code in the **JMDplots** package (<https://github.com/jedick/JMDplots>).
2. Development in 2022. Dick and Meng (2023) compared community  $Z_C$  with redox potential measurements from local to global scales. The term “community reference proteomes” was first applied, and **chem16S** was split into a separate package.
3. Late 2022. GTDB was added as a reference database.
4. June–July 2023. Integration with **phyloseq** and addition of vignettes: *Chemical metrics of reference proteomes*, *Integration of chem16S with phyloseq*, and *Plotting two chemical metrics*. Default reference database changed to GTDB.

### Options set in package chem16S

**chem16S** sets an option using the global `options` mechanism in R. This option will be set when package **chem16S** (or its namespace) is loaded if not already set.

`manual_mappings` A data frame of mappings between RDP and NCBI (RefSeq) taxonomies, which is read from `'extdata/manual_mappings.csv'`. The columns include `rdp.rank`, `rdp.name`, `ncbi.rank`, `ncbi.name`, and `notes`. This option is made available so the user can modify the manual mappings used by `map_taxa` at runtime.

### Files in extdata/RefSeq

NOTE: None of the `'*.R'` files in the `'extdata'` directories are included in the package submitted to CRAN; see GitHub or Zenodo for these files.

This directory contains two sets of files: 1) scripts to process source RefSeq sequence files to generate amino acid compositions of species-level reference proteomes and taxonomic names; 2) script and output for amino acid compositions of higher-level taxa. The files are based on RefSeq release 206 of 2021-05-21 (O'Leary et al., 2016).

**'README.txt'** Description of steps to generate reference proteomes of species-level taxa (including downloads and shell commands).

**'gencat.sh'** Helper script to extract microbial protein records from the RefSeq catalog.

**'protein\_refseq.R'** R code to sum the amino acid compositions of all proteins for each bacterial, archaeal, and viral species in the NCBI Reference Sequence database. NOTE: To save space, the output file (`'protein_refseq.csv'`) is not stored here but is in the `extdata/RefSeq` directory of the JMDplots package on GitHub (<https://github.com/jedick/JMDplots>). The first five columns are: `protein` ("refseq"), `organism` (taxonomic id), `ref` (organism name), `abbrv` (empty), `chains` (number of protein sequences for this organism). Columns 6 to 25 have the counts of amino acids.

**'taxon\_names.R'** R code for processing taxonomic IDs; the output file is `'taxon_names.csv'`. The columns are NCBI taxonomic ID (`taxid`), and names at different taxonomic rank (species, genus, family, order, class, phylum, superkingdom).

**'taxon\_AA.R'** Functions to create the files listed below:

**'taxon\_AA.csv.xz'** Average amino acid composition of reference proteomes for all species in each genus, family, order, class, phylum, and superkingdom.

**'taxon\_metrics.csv.xz'** Selected chemical metrics ( $Z_C$  and  $n_{H_2O}$ ) of reference proteomes for taxa at genus and higher ranks.

### Files in extdata/GTDB

**'taxon\_AA.R'** Functions to process GTDB source files (Parks et al., 2022) and produce the following output file:

**'taxon\_AA.csv.xz'** Average amino acid composition of reference proteomes for all species in each genus, family, order, class, phylum, and domain. In both this file and the corresponding file for RefSeq (see above), the `protein`, `organism`, `ref`, and `abbrv` columns contain the rank, taxon name, number of species used to generate the amino acid composition of this taxon, and parent taxon. `chains` is 1, denoting a single polypeptide chain, so the amino acid composition represents the average per-protein amino acid composition in this taxon, and the sum of amino acid counts is the average protein length.

**Files in extdata/metadata**

**'BGPF13.csv'** Metadata for Heart Lake Geyser Basin, Yellowstone (Bowen De León et al., 2012).

**'HLA+16.csv'** Metadata for the Baltic Sea (Herlemann et al., 2016).

**'SMS+12.csv'** Metadata for Bison Pool, Yellowstone (Swingley et al., 2012).

**Files in extdata/RDP**

Output of RDP Classifier with the default training set.

**'pipeline.R'** Pipeline for sequence data processing (uses external programs fastq-dump, vsearch, seqtk, RDP Classifier). This was used to make the files in both 'RDP' and 'RDP-GTDB' (the latter with GTDB = TRUE in the script).

**'BGPF13.tab.xz'** Heart Lake Geyser Basin.

**'HLA+16.tab.xz'** Baltic Sea.

**'SMS+12.tab.xz'** Bison Pool.

**Files in extdata/RDP-GTDB**

Output of RDP Classifier using the GTDB SSU training set ([doi:10.5281/zenodo.7633100](https://doi.org/10.5281/zenodo.7633100)).

**'BGPF13.tab.xz'** Heart Lake Geyser Basin.

**'HLA+16.tab.xz'** Baltic Sea.

**'SMS+12.tab.xz'** Bison Pool.

**Files in extdata/DADA2**

Identification and taxonomic classification of sequences using DADA2.

**'FEN+22'** Analysis of data from Fonseca et al. (2022) for marine sediment from the Humboldt Sulfuretum. 'pipeline.R' has the commands used to process the 16S rRNA gene sequence data and was adapted by Jeffrey Dick from the DADA2 pipeline tutorial (Callahan, 2020). 'SraRunInfo.csv' was obtained from the NCBI Sequence Read Archive (SRA) (<https://www.ncbi.nlm.nih.gov/sra/?term=PRJNA251688>). 'sample\_data.csv' has data obtained from NCBI BioSample records for BioProject PRJNA251688. '\*.png' are several plots created while running the DADA2 pipeline. 'ps\_FEN+22.rds' contains the phyloseq object with (including otu\_table, sample\_data, and refseq objects) created at the end of the DADA2 pipeline.

**'ZFZ+23'** Analysis of data from Zhang et al. (2023) for hot springs in the Qinghai-Tibet Plateau. 'pipeline.R' has the commands used to process the 16S rRNA gene sequence data and was adapted by Jeffrey Dick from the DADA2 pipeline tutorial (Callahan, 2020). 'SraRunInfo.csv' was obtained from the NCBI Sequence Read Archive (SRA) (<https://www.ncbi.nlm.nih.gov/sra/?term=PRJNA860942>). 'sample\_data.csv' has data obtained from NCBI BioSample records for BioProject PRJNA860942. '\*.png' are several plots created while running the DADA2 pipeline. 'ps\_ZFZ+23.rds' contains the phyloseq object with (including otu\_table, sample\_data, and refseq objects) created at the end of the DADA2 pipeline.

## References

- Bowen De León K, Gerlach R, Peyton BM, Fields MW. 2013. Archaeal and bacterial communities in three alkaline hot springs in Heart Lake Geysir Basin, Yellowstone National Park. *Frontiers in Microbiology* **4**: 330. doi:10.3389/fmicb.2013.00330
- Callahan B. 2020. DADA2 Pipeline Tutorial (1.16). <https://benjjneb.github.io/dada2/tutorial.html>, accessed on 2023-06-14.
- Dick JM, Meng D. 2023. Community- and genome-based evidence for a shaping influence of redox potential on bacterial protein evolution. *mSystems*. doi:10.1128/msystems.0001423
- Dick JM, Tan J. 2023. Chemical links between redox conditions and estimated community proteomes from 16S rRNA and reference protein sequences. *Microbial Ecology* **85**: 1338–1355. doi:10.1007/s00248022019889
- Fonseca A, Espinoza C, Nielsen LP, Marshall IPG, Gallardo VA. 2022. Bacterial community of sediments under the Eastern Boundary Current System shows high microdiversity and a latitudinal spatial pattern. *Frontiers in Microbiology* **13**: 1016418. doi:10.3389/fmicb.2022.1016418
- Herlemann DPR, Lundin D, Andersson AF, Labrenz M, Jürgens K. 2016. Phylogenetic signals of salinity and season in bacterial community composition across the salinity gradient of the Baltic Sea. *Frontiers in Microbiology* **7**: 1883. doi:10.3389/fmicb.2016.01883
- O’Leary NA et al. 2016. Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids Research* **44**: D733–D745. doi:10.1093/nar/gkv1189
- Parks DH, Chuvochina M, Rinke C, Mussig AJ, Chaumeil P-A, Hugenholtz P. 2022. GTDB: an ongoing census of bacterial and archaeal diversity through a phylogenetically consistent, rank normalized and complete genome-based taxonomy. *Nucleic Acids Research* **50**: D785–D794. doi:10.1093/nar/gkab776
- Swingley WD, Meyer-Dombard DR, Shock EL, Alsop EB, Falenski HD, Havig JR, Raymond J. 2012. Coordinating environmental genomics and geochemistry reveals metabolic transitions in a hot spring ecosystem. *PLOS One* **7**(6): e38108. doi:doi.org/10.1371/journal.pone.0038108
- Zhang H-S, Feng Q-D, Zhang D-Y, Zhu G-L, Yang L. 2023. Bacterial community structure in geothermal springs on the northern edge of Qinghai-Tibet plateau. *Frontiers in Microbiology* **13**: 994179. doi:10.3389/fmicb.2022.994179

---

calc\_metrics

*Calculate chemical metrics from amino acid composition*

---

## Description

Calculates various chemical metrics from amino acid composition(s) of one or more proteins.

## Usage

```
calc_metrics(AAcomp, metrics = c("Zc", "nO2", "nH2O"))
```

## Arguments

AComp	data frame of amino acid composition with column names being the 3-letter abbreviations of the amino acids (first letter uppercase)
metrics	character, chemical metrics to calculate

## Details

The following metrics are available:

- Zc: Carbon oxidation state. Its value is calculated from elemental composition as described by Dick (2014).
- nO2: Stoichiometric oxidation state. nO2 and nH2O are calculated from theoretical formation reactions of proteins from the QEC basis species (glutamine, glutamic acid, cysteine, H<sub>2</sub>O, and O<sub>2</sub>) and are normalized by number of amino acid residues (see Dick et al., 2020).
- nH2O: Stoichiometric hydration state. See above; note that the contribution of terminal groups is not counted in the calculation of nH2O.
- GRAVY: Grand average of hydropathicity. Values of the hydropathy index for individual amino acids are from Kyte and Doolittle (1982).
- pI: Isoelectric point. The net charge for each ionizable group was pre-calculated from pH 0 to 14 at intervals of 0.01. The isoelectric point is found as the pH where the sum of charges of all groups in the protein is closest to zero. The pK values for the terminal groups and sidechains are taken from Bjellqvist et al. (1993) and Bjellqvist et al. (1994); note that the calculation does not implement position-specific adjustments described in the latter paper. The number of N- and C-terminal groups is taken to be one, unless a column for chains (number of polypeptide chains) is given in AComp.
- MW: Molecular weight. This is the per-residue molecular weight and doesn't include the contribution of terminal groups.
- length or Length: Protein length. This is the number of amino acid residues per protein.
- H/C, H\_C, or HC: H/C ratio (not counting terminal H-OH groups of the protein backbone).
- N/C, N\_C, or NC: N/C ratio.
- O/C, O\_C, or OC: O/C ratio (not counting terminal H-OH groups of the protein backbone).
- S/C, S\_C, or SC: S/C ratio.

## Value

A data frame with the same number of rows as AComp and one column of numeric values for each of the metrics. An error is produced if any of the metrics is not available for calculation.

## References

- Bjellqvist B, Hughes GJ, Pasquali C, Paquet N, Ravier F, Sanchez J-C, Frutiger S, Hochstrasser D. 1993. The focusing positions of polypeptides in immobilized pH gradients can be predicted from their amino acid sequences. *Electrophoresis* **14**: 1023–1031. doi:10.1002/elps.11501401163
- Bjellqvist B, Basse B, Olsen E, Celis JE. 1994. Reference points for comparisons of two-dimensional maps of proteins from different human cell types defined in a pH scale where isoelectric points correlate with polypeptide compositions. *Electrophoresis* **15**: 529–539. doi:10.1002/elps.1150150171

Dick JM. 2014. Average oxidation state of carbon in proteins. *J. R. Soc. Interface* **11**: 20131095. doi:10.1098/rsif.2013.1095

Dick JM, Yu M, Tan J. 2020. Uncovering chemical signatures of salinity gradients through compositional analysis of protein sequences. *Biogeosciences* **17**: 6145–6162. doi:10.5194/bg1761452020

Gasteiger E, Hoogland C, Gattiker A, Duvaud S, Wilkins MR, Appel RD, Bairoch A. 2005. Protein identification and analysis tools on the ExPASy server. In: Walker JM, editor. *The Proteomics Protocols Handbook*. Totowa, NJ: Humana Press. pp. 571–607. doi:10.1385/1592598900:571

Kyte J, Doolittle RF. 1982. A simple method for displaying the hydropathic character of a protein. *J. Mol. Biol.* **157**: 105–132. doi:10.1016/00222836(82)905150

## See Also

[get\\_metrics](#)

## Examples

```
# Amino acid composition of alanylglycine
AG <- data.frame(Ala = 1, Gly = 1)
# Calculate default metrics
calc_metrics(AG)
# Calculate selected metrics
calc_metrics(AG, c("GRAVY", "pI"))
```

---

chem16S-data

*'phyloseq-class' objects generated using the DADA2 Pipeline Tutorial*

---

## Description

Objects in [phyloseq-class](#) format created by following the DADA2 Pipeline Tutorial (Callahan, 2020).

## Details

The example dataset for gut communities in a single mouse was taken from the [mothur MiSeq SOP](#). It is an extract of the complete dataset reported by Schloss et al. (2012). The output files were generated by using three different training sets to assign taxonomy to genus level:

**'mouse.silva'** silva\_nr99\_v138.1: [Silva 138.1 prokaryotic SSU taxonomic training data formatted for DADA2](#)

**'mouse.RDP'** rdp\_train\_set\_18: [RDP taxonomic training data formatted for DADA2 \(RDP trainset 18/release 11.5\)](#)

**'mouse.GTDB'** GTDB\_bac120\_arc53\_ssu\_r207: [GTDB: DADA2 formatted 16S rRNA gene sequences for both bacteria & archaea](#)

## Author(s)

Jeffrey M. Dick, following the DADA2 Pipeline Tutorial by Callahan (2020).

## References

Callahan B (2020). DADA2 Pipeline Tutorial (1.16). <https://benjjneb.github.io/dada2/tutorial.html>, accessed on 2023-06-14. Publication year taken from [https://code.bioconductor.org/browse/dada2/blob/RELEASE\\_3\\_11/DESCRIPTION](https://code.bioconductor.org/browse/dada2/blob/RELEASE_3_11/DESCRIPTION).

Schloss PD, Schubert AM, Zackular JP, Iverson KD, Young VB, Petrosino JF. 2012. Stabilization of the murine gut microbiome following weaning. *Gut Microbes* **3**(4): 383–393. doi:10.4161/gmic.21008

## Examples

```
# Load one data set
data(mouse.silva)
# Summarize the object (phyloseq-class)
mouse.silva
```

---

chemlab

*Formatted labels for chemical metrics*

---

## Description

Formatted labels for indicating chemical metrics on plots.

## Usage

```
chemlab(varname)
```

## Arguments

varname            character, the name of a variable

## Details

This function provides formatted labels for chemical metrics and related variables. varname can be any one of the chemical metrics described in [calc\\_metrics](#).

## Value

For most values of varname, a language object suitable for plotting (see [plotmath](#)). For labels where no special mathematical notation is required, a character object. If varname is not among the chemical metrics described in [calc\\_metrics](#), the value of varname itself.

## See Also

This function is used by [plot\\_metrics](#), [plot\\_ps\\_metrics](#), and [plot\\_ps\\_metrics2](#).

## Examples

```
chemlab("Zc")
```



---

get_metadata	<i>Template function for metadata</i>
--------------	---------------------------------------

---

## Description

Simple function that processes metadata files.

## Usage

```
get_metadata(file, metrics)
```

## Arguments

file	character, metadata file (CSV format)
metrics	data frame, output of <code>get_metrics</code>

## Details

This function is provided as a simple example of reading metadata about 16S rRNA datasets. The basic metadata for all datasets include 'name' (study name), 'Run' (SRA or other run ID), 'BioSample' (SRA BioSample), and 'Sample' (sample description). Each data set contains environmental data, but the specific variables differ between datasets. Therefore, this function processes the environmental data for different datasets and appends 'pch' and 'col' columns that can be used for visualization and identifying trends between groups of samples.

Samples are identified by the Run column (case-sensitive) in both file and metrics. It is possible that metrics for one or more samples are not available (for example, if the number of RDP classifications is below `mincount` in `get_metrics`). The function therefore removes metadata for unavailable samples. Then, the samples in metrics are placed in the same order as they appear in metadata. Ordering the samples is important because the columns in RDP Classifier output are not necessarily in the same order as the metadata.

## Value

If metrics is NULL, a data frame read from file with columns 'pch' and 'col' appended. Otherwise, a list with two components, metadata and metrics. In this list, metrics is the value from the metrics argument and metadata is a data frame read from file with columns 'pch' and 'col' appended and excluding any rows that correspond to unavailable samples in metrics.

## Examples

```
# Get metrics for the Bison Pool dataset (Swingley et al., 2012)
RDPfile <- system.file("extdata/RDP/SMS+12.tab.xz", package = "chem16S")
RDP <- read_RDP(RDPfile)
map <- map_taxa(RDP, refdb = "RefSeq")
metrics <- get_metrics(RDP, map, refdb = "RefSeq")
# Read the metadata file
mdatfile <- system.file("extdata/metadata/SMS+12.csv", package = "chem16S")
mdat <- get_metadata(mdatfile, metrics)
```

```

# Print sample names
mdat$metadata$Sample # 1 2 3 4 5
# Plot Zc of community reference proteomes at sites 1-5
Site <- mdat$metadata$Sample
Zc <- mdat$metrics$Zc
plot(Site, Zc, ylab = "Carbon oxidation state", type = "b")

# Find the columns of the RDP file for samples 1 to 5
match(mdat$metadata$Run, colnames(RDP)) # 6 5 7 8 9
# NOTE: for this dataset, samples 1 and 2 are in columns 6 and 5, respectively

```

---

get\_metrics

*Calculate chemical metrics of community reference proteomes*


---

## Description

Combines RDP classifications with RefSeq reference proteomes to get amino acid compositions of community reference proteomes. Amino acid compositions are used to calculate chemical metrics.

## Usage

```

get_metrics(RDP, map, refdb = "GTDB", taxon_AA = NULL, groups = NULL,
  return_AA = FALSE, zero_AA = NULL, metrics = c("Zc", "nO2", "nH2O"))

```

## Arguments

RDP	data frame, taxonomic abundances produced by <a href="#">read_RDP</a> or <a href="#">ps_taxaccounts</a>
map	data frame, taxonomic mapping produced by <a href="#">map_taxa</a>
refdb	character, name of reference database ('GTDB' or 'RefSeq')
taxon_AA	data frame, contents of 'taxon_AA.csv.xz' in 'extdata/GTDB' or 'extdata/RefSeq' (optionally used for speed-up)
groups	list of indexing vectors, samples to be aggregated into groups
return_AA	logical, return the amino acid composition for each sample instead of the chemical metrics?
zero_AA	character, three-letter abbreviation(s) of amino acid(s) to assign zero counts for calculating chemical metrics
metrics	character, the chemical metrics to calculate

## Details

`get_metrics` calculates selected chemical metrics for the community reference proteome in each sample. The community reference proteome is computed from the amino acid compositions of reference proteomes for taxa (obtained from the reference database in `refdb`), multiplied by taxonomic abundances given in RDP. RDP may include results from the RDP Classifier (read using [read\\_RDP](#)) or derived from the OTU table of a [phyloseq-class](#) object (see [ps\\_taxaccounts](#)). `map` defines

the taxonomic mapping between RDP and refdb. Then, chemical metrics are calculated from the amino acid composition of the community reference proteome. The default chemical metrics are carbon oxidation state ( $Z_C$ ), stoichiometric oxidation state ( $n_{O_2}$ ), and stoichiometric hydration state ( $n_{H_2O}$ ). See [calc\\_metrics](#) for other available metrics.

groups, if given, is a list of one or more indexing vectors (with logical or numeric values) corresponding to samples whose taxonomic classifications are aggregated into groups before calculating amino acid compositions and chemical metrics.

### Value

A data frame with one row for each sample, corresponding to columns 5 and above of RDP. The sample names are in the first column, which is named Run by default, or group if the groups argument is provided. The remaining columns have numeric values and are named for each of the calculated metrics.

### See Also

[calc\\_metrics](#), [plot\\_metrics](#)

### Examples

```
## First two examples are for RDP Classifier with default training set
## and mapping to NCBI taxonomy with RefSeq reference proteomes

# Get chemical metrics for all samples in a dataset
RDPfile <- system.file("extdata/RDP/BGPF13.tab.xz", package = "chem16S")
RDP <- read_RDP(RDPfile)
map <- map_taxa(RDP, refdb = "RefSeq")
# This is a data frame with 14 rows and Run, Zc, nO2, and nH2O columns
(metrics <- get_metrics(RDP, map, refdb = "RefSeq"))

# Read the metadata file
mdatfile <- system.file("extdata/metadata/BGPF13.csv", package = "chem16S")
# Create list with metadata and metrics in same sample order
mdat <- get_metadata(mdatfile, metrics)
# Calculate metrics for aggregated samples of Archaea and Bacteria
groups <- list(A = mdat$metadata$domain == "Archaea",
  B = mdat$metadata$domain == "Bacteria")
# This is a data frame with 2 rows and group, Zc, nO2, and nH2O columns
get_metrics(RDP, map, refdb = "RefSeq", groups = groups)

## For the next example, classifications were made using the RDP Classifier
## with the GTDB SSU training set, and we use reference proteomes based on GTDB
## (this uses the default option of refdb = "GTDB").

RDPfile.GTDB <- system.file("extdata/RDP-GTDB/BGPF13.tab.xz", package = "chem16S")
RDP.GTDB <- read_RDP(RDPfile.GTDB)
map.GTDB <- map_taxa(RDP.GTDB)
metrics.GTDB <- get_metrics(RDP.GTDB, map.GTDB)

# Plot Zc from GTDB vs RefSeq
```

```

xylim <- range(metrics$Zc, metrics.GTDB$Zc)
plot(metrics$Zc, metrics.GTDB$Zc, xlim = xylim, ylim = xylim, type = "n")
lines(xylim, xylim, lty = 2, col = 8)
points(metrics$Zc, metrics.GTDB$Zc, pch = mdat$metadata$pch, bg = mdat$metadata$col)
md.leg <- mdat$metadata[1:2, ]
legend("bottomright", md.leg$domain, pch = md.leg$pch, pt.bg = md.leg$col)
title(quote(italic(Z)[C]~"from GTDB vs RefSeq"))

# To exclude tryptophan, tyrosine, and phenylalanine
# from the calculation of chemical metrics
get_metrics(RDP.GTDB, map.GTDB, zero_AA = c("Trp", "Tyr", "Phe"))

```

map\_taxa

*Map taxonomic names to RefSeq or GTDB taxonomy***Description**

Maps taxonomic names to RefSeq (NCBI) or GTDB taxonomy by automatic matching of taxonomic names, with manual mappings for some groups.

**Usage**

```
map_taxa(taxacounts = NULL, refdb = "GTDB", quiet = FALSE)
```

**Arguments**

taxacounts	data frame with taxonomic name and abundances
refdb	character, name of reference database ('GTDB' or 'RefSeq')
quiet	logical, suppress printed messages?

**Details**

This function maps taxonomic names to the RefSeq (NCBI) or GTDB taxonomy. taxacounts should be a data frame generated by either [read\\_RDP](#) or [ps\\_taxacounts](#). Input names are made by combining the taxonomic rank and name with an underscore separator (e.g. 'genus\_Escherichia/Shigella'). Input names are then matched to the taxa listed in 'taxon\_AA.csv.xz' found under 'extdata/RefSeq' or 'extdata/GTDB'. The protein and organism columns in these files hold the rank and taxon name extracted from the RefSeq or GTDB database. Only exactly matching names are automatically mapped.

For mapping to the RefSeq (NCBI) taxonomy, some group names are manually mapped as follows (see Dick and Tan, 2023):

**Input (i.e., RDP)**

genus\_Escherichia/Shigella  
 phylum\_Cyanobacteria/Chloroplast  
 genus\_Marinimicrobia\_genera\_incertae\_sedis  
 class\_Cyanobacteria

**RefSeq**

genus\_Escherichia  
 phylum\_Cyanobacteria  
 species\_Candidatus\_Marinimicrobia\_bacterium  
 phylum\_Cyanobacteria

genus_Spartobacteria_genera_incertae_sedis	species_Spartobacteria bacterium LR76
class_Planctomycetacia	class_Planctomycetia
class_Actinobacteria	phylum_Actinobacteria
order_Rhizobiales	order_Hyphomicrobiales
genus_Gp1	genus_Acidobacterium
genus_Gp6	genus_Luteitalea
genus_GpI	genus_Nostoc
genus_GpIIa	genus_Synechococcus
genus_GpVI	genus_Pseudanabaena
family_Family II	family_Synechococcaceae
genus_Subdivision3_genera_incertae_sedis	family_Verrucomicrobia subdivision 3
order_Clostridiales	order_Eubacteriales
family_Ruminococcaceae	family_Oscillospiraceae

To avoid manual mapping, GTDB can be used for both taxonomic assignments and reference proteomes. Taxonomic assignments can be made using the RDP Classifier with this GTDB SSU training set: [doi:10.5281/zenodo.7633100](https://doi.org/10.5281/zenodo.7633100) or **dada2** with this GTDB training set: [doi:10.5281/zenodo.6655692](https://doi.org/10.5281/zenodo.6655692). Example files created using the RDP Classifier are provided under 'extdata/RDP-GTDB'. An example dataset created with DADA2 is `data(mouse.GTDB)`; this is a `phyloseq-class` object that can be processed with functions described at [physeq](#).

Change `quiet` to `TRUE` to suppress printing of messages about manual mappings, most abundant unmapped groups, and overall percentage of mapped names.

## Value

Integer vector with length equal to number of rows of taxacounts. Values are rownumbers in the data frame generated by reading `taxon_AA.csv.xz`, or `NA` for no matching taxon. Attributes `unmapped_groups` and `unmapped_percent` have the input names of unmapped groups and their percentage of the total classification count.

## References

Dick JM, Tan J. 2023. Chemical links between redox conditions and estimated community proteomes from 16S rRNA and reference protein sequences. *Microbial Ecology* **85**: 1338–1355. [doi:10.1007/s00248022019889](https://doi.org/10.1007/s00248022019889)

## Examples

```
# Mapping taxonomic classifications from RDP training set to NCBI taxonomy
file <- system.file("extdata/RDP/SMS+12.tab.xz", package = "chem16S")
RDP <- read_RDP(file)
map <- map_taxa(RDP, refdb = "RefSeq")
# About 24% of classifications are unmapped
sum(attributes(map)$unmapped_percent)

# Mapping from GTDB training set to GTDB taxonomy
file <- system.file("extdata/RDP-GTDB/SMS+12.tab.xz", package = "chem16S")
RDP.GTDB <- read_RDP(file)
map.GTDB <- map_taxa(RDP.GTDB)
```

```
# There is 100% mapping (zero unmapped classifications)
sum(attributes(map.GTDB)$unmapped_percent)
```

---

physeq

*Process 'phyloseq-class' objects to calculate chemical metrics*


---

## Description

Calculate chemical metrics of community reference proteomes from taxonomic abundances in [phyloseq-class](#) objects and make plots of individual metrics or two metrics against each other.

## Usage

```
ps_taxacounts(physeq, split = TRUE)
ps_metrics(physeq, split = TRUE, refdb = "GTDB", quiet = FALSE, ...)
plot_ps_metrics(physeq, x = "samples", color = NULL, shape = NULL,
  title = NULL, scales = "free_y", nrow = 1,
  metrics = c("Zc", "nO2", "nH2O"), sortby = NULL, refdb = "GTDB",
  quiet = FALSE)
plot_ps_metrics2(physeq, x = "Zc", y = "nH2O", color = NULL,
  shape = NULL, title = NULL, refdb = "GTDB", quiet = FALSE)
```

## Arguments

physeq	'phyloseq-class', a <b>phyloseq</b> object containing an 'otu_table' and 'tax_table'
split	logical, calculate metrics for each sample (TRUE) or the entire set (FALSE)
refdb	character, name of reference database ('GTDB' or 'RefSeq'; passed to <a href="#">get_metrics</a> )
quiet	logical, suppress printed messages? (passed to <a href="#">map_taxa</a> )
...	additional arguments for <a href="#">get_metrics</a>
x	map this variable to the horizontal axis
color	map this variable to colors
shape	map this variable to shapes
title	main title for the plot
scales	'free_y' or 'fixed' for different vertical scales or same scale in each panel
nrow	number of rows of panels (passed to <a href="#">facet_wrap</a> )
metrics	character, chemical metrics to calculate
sortby	subset of metrics argument, sort x-index by the mean values of these metrics
y	map this variable to the vertical axis

## Details

`ps_taxaccounts` returns a data frame with lowest-level classifications (from genus to phylum) and abundances for each OTU in the `physeq` object.

`ps_metrics` calculates chemical metrics of community reference proteomes from the lowest-level classifications after taxonomic mapping with `map_taxa`. To choose the metrics to calculate, use the `metrics` argument, which is passed via ... to `get_metrics`.

`plot_ps_metrics` plots individual chemical metrics or multiple chemical metrics, each in their own panel. This function is modified from `plot_richness`; see that help page for further details about graphics-related arguments.

`plot_ps_metrics2` plot two chemical metrics against each other.

## Value

For `ps_taxaccounts`, a data frame with samples in rows, and columns `taxid` (sample name taken from the OTU table obtained from `otu_table(physeq)`), `lineage` (all taxonomic names from the taxonomy table obtained from `tax_table(physeq)`, separated by semicolons), `name` (lowest-level taxonomic name), `rank` (rank of lowest-level taxon), followed by the abundances from the OTU table.

For `ps_metrics`, a data frame with samples in rows, and columns corresponding to each of the metrics.

For `plot_ps_metrics` and `plot_ps_metrics2`, ggplot objects.

## Author(s)

The author of this Rd file and the `ps_taxaccounts` and `ps_metrics` functions is Jeffrey Dick. `plot_ps_metrics` is based on the `plot_richness` function by Paul J. McMurdie and was modified from that function by Jeffrey Dick.

## Examples

```
# Example 1: Humboldt Sulfuretum (Fonseca et al., 2022)
file <- system.file("extdata/DADA2/FEN+22/ps_FEN+22.rds", package = "chem16S")
physeq <- readRDS(file)
# Get lowest-level (to genus) classification for each OTU
taxaccounts <- ps_taxaccounts(physeq)
# Show numbers of assignments at each taxonomic level
table(taxaccounts$rank)
# Map taxonomic names
map <- map_taxa(taxaccounts)
# Taxonomy and reference proteomes are both from GTDB,
# so there are no unmapped classifications
sum(attr(map, "unmapped_percent")) # == 0
# Calculate chemical metrics
ps_metrics(physeq)

# Example 2: GlobalPatterns dataset from phyloseq
data(GlobalPatterns, package = "phyloseq")
# Plot metrics grouped by sample type and sorted by mean Zc;
# refdb = "RefSeq" uses manual mappings from the RDP to NCBI taxonomy
```

```
p <- plot_ps_metrics(GlobalPatterns, x = "SampleType",
                    sortby = "Zc", refdb = "RefSeq")
# Change orientation of x-axis labels
p + ggplot2::theme(axis.text.x = ggplot2::element_text(
  angle = 45, vjust = 1, hjust = 1))
```

---

plot\_metrics

*Plot chemical metrics of community reference proteomes*


---

## Description

Functions to plot chemical metrics of community reference proteomes.

## Usage

```
plot_metrics(mdat, identify = FALSE, title = TRUE,
            xlim = NULL, ylim = NULL, plot.it = TRUE, points = TRUE, lines = FALSE,
            cex = 1, pch1 = 1, pch2 = 21, return = "data",
            extracolumn = NULL, add = FALSE, pt.open.col = 1,
            xlab = chemlab("Zc"), ylab = chemlab("nH2O"))
add_hull(x, y, basecol, outline = FALSE, ...)
```

## Arguments

mdat	list, output by <a href="#">get_metadata</a>
identify	logical, run <a href="#">identify</a> for interactive identification of points?
title	character, plot title
xlim	numeric, x axis limits
ylim	numeric, y axis limits
plot.it	logical, make a plot?
points	logical, plot points?
lines	logical, plot lines?
cex	numeric, point size
pch1	numeric, symbol for samples in group 1
pch2	numeric, symbol for samples in group 2
return	character, indicates whether to return 'data' values or group 'means'
extracolumn	character, the name of one or more extra columns (from <a href="#">get_metadata</a> ) to include in the output
add	logical, add to existing plot?
pt.open.col	color of border for open point symbols (pch > 20)
xlab	x axis label
ylab	y axis label



x	x values
y	y values
basecol	color used for hull lines
outline	logical, draw hull lines?
...	arguments passed to <a href="#">polygon</a>

### Details

plot\_metrics plots the  $Z_C$  and  $n_{H_2O}$  values provided in `mdat$metrics`. Symbol shape and color (`pch` and `col`) are taken from `mdat$metadata`.

If `pch1` and `pch2` are provided, then samples are classified into two groups according to the value of `mdat$metadata$pch`. Mean values of the chemical metrics for each group are plotted with star-shaped symbols.

`add_hull` uses [polygon](#) to draw a convex hull around the points given in `x`, `y`. The lines are drawn only if `outline` is TRUE; otherwise, the area of the hull is filled. The fill color is a semi-transparent version of the color given in `basecol`.

### Value

For `plot_metrics`, a data frame with columns with study name and Run IDs ('name', 'Run'), plotting values ('nH2O', 'Zc'), and plotting symbols and colors ('pch', 'col'). For `add_hull`, no return value (called for side effects).

### References

Herlemann, D. P. R., Lundin, D., Andersson, A. F., Labrenz, M. and Jürgens, K. (2016) Phylogenetic signals of salinity and season in bacterial community composition across the salinity gradient of the Baltic Sea. *Front. Microbiol.* **7**, 1883. doi:10.3389/fmicb.2016.01883

### Examples

```
# Make a plot for the Baltic Sea salinity gradient
# (data from Herlemann et al., 2016)
RDPfile <- system.file("extdata/RDP/HLA+16.tab.xz", package = "chem16S")
RDP <- read_RDP(RDPfile)
map <- map_taxa(RDP, refdb = "RefSeq")
metrics <- get_metrics(RDP, map, refdb = "RefSeq")
mdatfile <- system.file("extdata/metadata/HLA+16.csv", package = "chem16S")
mdat <- get_metadata(mdatfile, metrics)
pm <- plot_metrics(mdat)
# Add a legend
legend <- c("< 6 PSU", "6-20 PSU", "> 20 PSU")
pch <- c(24, 20, 21)
pt.bg <- c(3, NA, 4)
legend("bottomright", legend, pch = pch, col = 1, pt.bg = pt.bg, bg = "white")
# Classify samples with low and high salinity
ilo <- mdat$metadata$salinity < 6
ihi <- mdat$metadata$salinity > 20
# Add convex hulls
```

```

add_hull(pm$Zc[i1o], pm$nH2O[i1o], "green3")
add_hull(pm$Zc[ihi], pm$nH2O[ihi], "blue")

# Show points for all samples and larger star-shaped points
# for mean values of high- and low-salinity samples
plot_metrics(mdat, pch1 = 21, pch2 = 24)

# Make a plot for only Proteobacteria
RDP <- read_RDP(RDPfile, lineage = "Proteobacteria")
map <- map_taxa(RDP, refdb = "RefSeq")
metrics <- get_metrics(RDP, map, refdb = "RefSeq")
mdatfile <- system.file("extdata/metadata/HLA+16.csv", package = "chem16S")
mdat <- get_metadata(mdatfile, metrics)
mdat$metadata$name <- paste(mdat$metadata$name, "(Proteobacteria)")
plot_metrics(mdat)

```

---

read\_RDP

*Read and filter RDP Classifier output*


---

## Description

Reads RDP Classifier output and optionally extracts lineages, removes classifications below a certain count, or truncates classifications to lowest taxonomic level.

## Usage

```

read_RDP(file, lineage = NULL, mincount = 200,
         lowest.level = NULL, quiet = FALSE)

```

## Arguments

file	character, file name
lineage	character, regular expression for filtering lineages
mincount	integer, samples with less than this number of RDP classifications are excluded
lowest.level	character, truncate classifications at this taxonomic level
quiet	logical, suppress printed messages?

## Details

read\_RDP reads and filters RDP results for all samples in a study. The input file should be created by RDP Classifier using the -h option to create a hierarchical listing. Classifications for multiple samples can be combined into a single file using the merge-count command of RDP Classifier.

Only rows (lineages) with count > 0 for at least one sample are retained. Sequences with classification results at only the root or domain level (Root, Archaea, or Bacteria) are omitted because they provide poor taxonomic resolution. Sequences classified to the class Chloroplast or genera Chlorophyta or Bacillariophyta are also omitted because they have little correspondence with the NCBI taxonomy. Then, only columns (samples) with classification count >= mincount are retained. All

remaining sequences (those classified to genus or higher levels) can be used for mapping to the NCBI taxonomy.

The lineage text of the RDP Classifier looks like “Root;rootrank;Archaea;domain;Diapherotrites; phylum;Diapherotrites Incertae Sedis AR10;genus;”, so you can use `lineage = "Archaea"` to select the archaeal classifications or `lineage = "genus"` to select genus-level classifications.

Use the `lowest.level` argument to truncate the classifications to a level higher than genus. This argument does not reduce the number of classifications, but only trims the RDP lineages to the specified level. This may create duplicate lineages, for which the classification counts are summed, and only unique lineages are present in the returned data frame.

Change `quiet` to `TRUE` to suppress printing of messages about percentage classification to genus level, omitted sequences, and final range of total counts among all samples.

### Value

Data frame with columns inherited from `file` (from RDP output files: `taxid`, `rank`, `lineage`, `name`, and one column of classification counts for each sample). If any sample has less than `mincount` total counts, that column is omitted from the result. The rows in the result are subset from those in `file` by filtering operations.

NB: taxids in RDP files are not NCBI taxids.

### Examples

```
# An example file
file <- system.file("extdata/RDP/SMS+12.tab.xz", package = "chem16S")

# Default settings
r0 <- read_RDP(file)

# Suppress messages
r1 <- read_RDP(file, quiet = TRUE)

# Increase minimum count threshold
r2 <- read_RDP(file, mincount = 500)
# This should be 3 (i.e., 3 samples have less than 500 counts)
ncol(r1) - ncol(r2)

# Keep only Archaea
r3 <- read_RDP(file, lineage = "Archaea")
# This should be TRUE
all(grepl("Archaea", r3$lineage))

# Truncate taxonomy at phylum
r4 <- read_RDP(file, lowest.level = "phylum")
# This should be TRUE
all(sapply(strsplit(r4$lineage, ";"), tail, 1) == "phylum")
```

# Index

## \* data

chem16S-data, 7

add\_hull (plot\_metrics), 16

calc\_metrics, 5, 8, 11

chem16S (chem16S-package), 2

chem16S-data, 7

chem16S-package, 2

chemlab, 8

facet\_wrap, 14

get\_metadata, 2, 9, 16

get\_metrics, 2, 7, 9, 10, 14, 15

identify, 16

map\_taxa, 2, 3, 10, 12, 14, 15

mouse.GTDB, 13

mouse.GTDB (chem16S-data), 7

mouse.RDP (chem16S-data), 7

mouse.silva (chem16S-data), 7

options, 3

otu\_table, 15

physeq, 2, 13, 14

plot\_metrics, 2, 8, 11, 16

plot\_ps\_metrics, 8

plot\_ps\_metrics (physeq), 14

plot\_ps\_metrics2, 8

plot\_ps\_metrics2 (physeq), 14

plot\_richness, 15

plotmath, 8

polygon, 17

ps\_metrics (physeq), 14

ps\_taxacounts, 10, 12

ps\_taxacounts (physeq), 14

read\_RDP, 2, 10, 12, 18

tax\_table, 15